

© 2010 Siddhartha Verma

DEVELOPMENT OF ERROR CORRECTION TECHNIQUES FOR NITRATE-N  
LOAD ESTIMATION MODELS

BY

SIDDHARTHA VERMA

THESIS

Submitted in partial fulfillment of the requirements  
for the degree of Master of Science in Agricultural and Biological Engineering  
in the Graduate College of the  
University of Illinois at Urbana-Champaign, 2010

Urbana, Illinois

Adviser:

Associate Professor Richard Cooke

# ABSTRACT

Excess nutrient loads carried by streams and rivers are a great concern for environmental resource managers. In agricultural regions, excess loads are transported downstream to receiving water bodies, potentially causing algal blooms, which could lead to numerous ecological problems. To better understand nutrient load transport, and to develop appropriate water management plans, it is important to have accurate estimates of annual nutrient loads.

This study used a Monte Carlo sub-sampling method and error-corrected statistical models to estimate annual nitrate-N loads from two watersheds in central Illinois. The performance of three load estimation methods (the seven-parameter log-linear model, the ratio estimator, and the flow-weighted averaging estimator) applied at one-, two-, four-, six-, and eight-week sampling frequencies were compared. Five error correction techniques; the existing composite method, and four new error correction techniques developed in this study; were applied to each combination of sampling frequency and load estimation method. On average, the most accurate error reduction technique, (proportional rectangular) resulted in 15% and 30% more accurate load estimates when compared to the most accurate uncorrected load estimation method (ratio estimator) for the two watersheds. Using error correction methods, it is possible to design more cost-effective monitoring plans by achieving the same load estimation accuracy with fewer observations.

Finally, the optimum combinations of monitoring threshold and sampling frequency that minimizes the number of samples required to achieve specified levels of accuracy in load estimation were determined. For one- to three-weeks sampling frequencies, combined threshold/fixed-interval monitoring approaches produced the best outcomes, while fixed-interval-only approaches produced the most accurate results for four- to eight-weeks sampling frequencies.

# ACKNOWLEDGEMENTS

I would like to thank everyone who has offered help, words of encouragement and support over the past two years thereby helping me complete this project.

Firstly, I would like to express my gratitude to my advisor Dr. Richard Cooke, who has been a constant source of motivation. His passion towards research, optimism and the ever smiling persona was highly inspirational for me over the past two years. He also provided me great opportunities which have helped me immensely to grow both personally and professionally.

I'll also like to thank Dr. Momcilo Markus for giving me this great opportunity to work on this project. He has been a wonderful mentor who always allowed me to share different ideas and approaches to research which has helped me become a better researcher. He was instrumental in keeping the project on track and ensuring its overall quality.

I'll also like to thank my committee members Dr. Prasanta Kalita and Dr. Luis Rodriguez, for taking time out to review my thesis and also for their invaluable suggestions and comments.

I am grateful to Dr. Amit Desai and Rabin Bhattarai who have been role models for me and also for addressing my innumerable technical and professional queries. I would also like to thank Dr. Paul Davidson, Greg Goodwin, Dan Koch, Dr. Huzefa Raja, Ajith Harish, Goutam Nistala, Steve Anderson and Greg Byard for sharing their graduate experiences with me. I am really thankful to Neha Bisht, who has always lent a patient ear and has helped me keep my priorities in line. Her enthusiasm and optimism have kept my spirits high throughout the past two years.

Finally I'll thank my parents and brother, without whom my journey thus far wouldn't have been possible. I'm really fortunate to have received their unwavering love and support throughout. They have always taught me to dream big and to work diligently to turn those dreams into reality.

# TABLE OF CONTENTS

Chapter 1	Introduction .....	1
Chapter 2	Review of literature .....	3
2.1	Crops .....	3
2.2	Fertilizers .....	4
2.3	Nutrient contaminants .....	5
2.4	Statistical Load Estimation Models.....	6
2.5	Composite Method .....	11
Chapter 3	Objectives.....	13
Chapter 4	Site and Dataset descriptions .....	14
4.1	Watershed Descriptions.....	14
4.2	Dataset Description .....	14
4.3	Flashiness Index .....	22
Chapter 5	Methodology .....	23
5.1	Load Estimation Methods .....	24
5.2	Autocorrelation of Modeling Residuals .....	27
5.3	Composite Method .....	29
5.4	Development of Error Correction Methods .....	30
5.5	Monte Carlo Simulation.....	34
5.6	Evaluation Criteria .....	35
Chapter 6	Results and Discussion.....	39
6.1	Load Calculations.....	39
6.2	Performance of Load Estimation Methods.....	39
6.3	Performance of Error Correction Techniques .....	44
6.4	Threshold Analysis.....	48
Chapter 7	Summary .....	52
Chapter 8	Conclusions .....	54
References	.....	55
Appendix A.	.....	61

A.1. Matlab code to calculate bias and RMSE for various load estimation models and error correction techniques using 700 iterations..... 61

A.2. Matlab code to compute load estimates from various load estimation models and error correction techniques..... 64

# CHAPTER 1

## INTRODUCTION

Increased nutrient application rates, coupled with intensive cropping patterns in the Midwestern United States, have been identified as important causes of hypoxia in the Gulf of Mexico. Substantial increases in riverine concentrations of nutrients such as nitrates and phosphorus have notably contributed to the hypoxia problem (Rabalais et al., 1996). Much of the agricultural lands in the Mississippi River watershed are drained with sub-surface (tile) drain systems designed to remove excess water quickly from agricultural fields. Nitrogen in the nitrate form is readily soluble in water, draining easily through tile drain systems (David et al., 1997). This nutrient-rich water, draining from throughout the Midwest, is discharged into agricultural ditches, then to small rivers, and is eventually carried by the Mississippi River into the Gulf of Mexico. Increased levels of nutrients in the Gulf of Mexico have led to eutrophication followed by hypoxia, eventually creating severe environmental, ecological, and economic imbalances (Turner and Rabalais, 1994).

A necessary step in the development of management practices to reduce nutrient transport from agricultural fields is an accurate characterization of nutrient loads in streams and rivers. However, there is a large degree of uncertainty involved with accurately estimating loads for any given waterway.

Nutrient load is primarily obtained by calculating the product of flow rate and solute concentration. Ideally, to attain the most accurate load estimates for a specific time interval for any stream or river, continuous datasets of flow rates and solute concentrations are needed. With advances in technology, robust sensors have been developed to record continuous flow rates in rivers and streams. On the other hand, economic considerations typically limit the frequency with which nutrient concentrations can be obtained. Researchers have developed various techniques to use sparse datasets of measured concentrations to obtain datasets that match the frequency of flow measurements. These techniques can be broadly classified into three major categories: (1) the use of deterministic models (e.g., Soil and Water Assessment Tool (SWAT), Hydrological Simulation Program-Fortran(HSPF)) (Arnold et al., 1998; Santhi et al., 2001;

Schilling and Wolter, 2009; Im et al., 2003); (2) the use of statistical models (e.g., seven-parameter regression model, ratio) (Preston et al., 1989; Mukhopadhyay and Smith, 2000; Guo et al., 2002; Quilbé et al., 2006) and (3) the use of Artificial Neural Networks (ANN) (Yu et al., 2004; Anctil et al., 2009). This research focuses on the use of statistical models for estimating nutrient loads.

In this study, several error-corrected statistical methods were developed and applied to the estimation of annual nitrate-N loads from two watersheds in Central Illinois, the Upper Sangamon River at Monticello (1,406 square kilometers [ $\text{km}^2$ ]) and the Vermilion River at Pontiac (1,570  $\text{km}^2$ ). The Monte Carlo sub-sampling technique was used to simulate weekly to bimonthly sampling frequencies for monitoring durations ranging from one to six years. For each unique combination of the monitoring duration and sampling frequency, three statistical models were compared; the seven-parameter log-linear model, the ratio estimator, and the flow-weighted averaging estimator. In addition, for each combination of monitoring duration, sampling frequency and statistical models, five error correction techniques were applied to find the most accurate load estimation procedure. In addition to the existing composite method, four new error correction techniques were developed; namely, the triangular-proportional, rectangular-proportional, triangular-residual, and rectangular-residual methods. Finally, this study examined the use of monitoring thresholds, above which discharges are continuously (daily) monitored. The purpose was to find optimum combinations of monitoring thresholds and sampling frequencies for which the load estimation errors were smallest.



# CHAPTER 2

## REVIEW OF LITERATURE

### 2.1 Crops

The extensive development of railroads in the early 20<sup>th</sup> century transformed Illinois agriculture from being a local scale supplier to a much wider national and international mass supplier. This transformation also brought a more industrialized outlook to farming and led to less diversification in agricultural production (Byard, 2009) with an increased focus on corn and soybean that emerged as the main agricultural products from the state of Illinois (Bramstedt and Endres, 1999).

Borah et al. (2003) reported agriculture to be the principal land use Illinois, with most of the land being used for agriculture in Illinois is used to grow corn and soybean. Typically, within the state, farmers follow a rotation pattern between corn and soybean over consecutive years. This corn-soybean rotation offers more flexibility in tillage and reduces yield losses due to large corn residues, as compared to continuous corn (Randall et al., 2002). Additionally, planting soybean also enriches the soil with nitrogen which can be attributed to the crop's symbiotic relationship with nitrogen fixing bacteria.

Corn is the most important agricultural product from Illinois, the second largest producer of corn in the United States. Typically, throughout the central part of the state corn is planted from the third week of April through May. Usually corn planting is accompanied by simultaneous herbicide application, but sometimes corn herbicides are also applied as pre plant applications 1 to 30 days prior to planting. Post-emergence applications are typically done after 4 to 6 weeks of planting (Byard, 2009).

Soybean production started in Illinois at the beginning of the 20<sup>th</sup> century. Over the past century, there has been a remarkable transformation in soybean production in the state and currently Illinois is the second largest producer of soybean in the United States. This notable increase in soybean production in the state can be credited to a suitable climate for its

production, coupled with a widespread infrastructure of ancillary industries for soy crushing and processing. As a large number of farmers in Illinois practice corn-soybean rotation, a considerable percentage of cropland in the state is used to grow soybean. Soybean planting is done at a later time as compared to corn planting, usually mid to late May, and nitrogen fertilizers are not normally applied (David et al., 2008).

Apart from corn and soybeans, farmers in the state of Illinois also grow and raise other agricultural commodities such as cattle, swine, wheat, vegetables, oats, and sorghum (Illinois Department of Agriculture, 2001).

## **2.2 Fertilizers**

The increased usage of fertilizers over the past few decades has been instrumental in significantly improving crop yields throughout the world. Nitrogen, phosphorus and potassium are the main macronutrients which are added to the soils as fertilizers. The majority of anthropogenic nitrogen and phosphorus in the environment is a result of extensive fertilizer application (Byard, 2009).

A majority of nitrogen application in Midwestern agricultural watersheds (temperate zones) takes place in the fall season. Cookson et al. (2000) reported that nitrate application in temperate areas around the world mainly boosts crop growth, decreases establishment times and also increases seed yields of crops. Nitrogen application as fertilizers can be done in various forms such as ammonia, nitrates and animal manure (Byard, 2009). Commercial nitrogen fertilizers are applied mainly in the form of ammonium and nitrates; with ammonium being readily convert into nitrates. Animal manure is also used as a source of nitrogen, but mainly in Southeast United States, whereas throughout the Midwest and Western United States commercial fertilizers dominate as sources of nitrogen.

In the United States, the sale of nitrogen fertilizers has been increasing steadily with a twenty fold increase between 1945 and 1985. Annually, approximately 11 million tons of commercial nitrogen fertilizer and 6.5 million tons of manure-based nitrogen fertilizer are applied in United States (USGS, 1996). Nitrate fertilizer application rates for Midwestern corn production varied from 150 to 225 kg per hectare per year (David et al., 2008).

Phosphorus is mainly applied as compounds of phosphate, and although it is not being very mobile, it is transported in significant quantities to streams by soil erosion. The use of phosphate fertilizers has increase fourfold between 1945 and 1985; with approximately 2 million tons being applied each year as commercial and manure-based phosphate fertilizer in United States (USGS, 1996).

## **2.3 Nutrient Contaminants**

The major sources of nutrients entering the environment can be classified as either point or non-point sources. Natural precipitation, dissolution of natural minerals from soil and geological formations, and fertilizer application are the main non-point sources credited with nutrient delivery to the environment. Of these sources of nutrient delivery to the environment, human influence is mainly seen in fertilizer application for agricultural production. Background nitrate concentrations are less than 2 mg/L for ground water, and less than 0.6 mg/L for streams, whereas the background total phosphorus concentrations are usually less than 0.1mg/L. (USGS, 1996). Fertilizer application mainly occurs in the form of nitrates and phosphates and any increase in the natural background concentrations of nitrates and total phosphorus is often classified as nutrient contamination

Over the many past decades, extensive development of subsurface (tile) drainage systems has taken place throughout the Midwestern United States. In Illinois alone, approximately 4 million hectares of land are artificially drained by tile drainage systems. These systems were designed to convert swamplands into highly productive farmlands, and although these systems promote increased yields, and reduced sediment losses, they increase the delivery of nitrates to receiving water bodies. Typically, subsurface drainage systems in Midwestern United States drain into a small agricultural ditches, which eventually drain into a river which in turn enters the Mississippi river (Kalita et al., 2007).

Nitrates are the most prevalent contaminant in groundwater in the United States, and, as they are readily soluble in water, they are transported in both surface and groundwater (USGS, 1996). Thus, discharge from the extensive subsurface drainage systems throughout the Midwest contains a large amount of nutrients which are eventually discharged into the Mississippi river system.

Both of the watersheds used in this study lie in the large Mississippi- Atchafalaya river system. The Mississippi-Atchafalaya river system drains 41% of conterminous United States. In the central part of this large river basin, i.e. the Midwest, most of the corn and soybean production of the country takes place and consequently a majority of fertilizers and pesticides used in the United States are applied within this river basin (Goolsby et al., 2001). Consequently, the Mississippi- Atchafalaya river system is the biggest source of riverine nutrient concentrations to the Gulf of Mexico.

The total nitrogen flux to the Gulf of Mexico for a period from 1980-1996 was reported to be 1,568,000 tons/yr. Out of the total nitrogen flux approximately 61% was in the nitrate form, 37% in the form of dissolved and particulate organic nitrogen and only 2% in the form of ammonium. Such elevated riverine concentrations of nitrates have been attributed as the one of the major causes of hypoxia in the Gulf of Mexico (Goolsby et al., 2001).

Although the role of nitrates in hypoxia is well documented, recently studies have indicated that phosphorus also plays a more important role than previously believed in causing hypoxia in the Gulf of Mexico (Sylvan et al., 2006; Alexander et al., 2008). Alexander et al. (2008) reported that analysis of different sources of nutrients to the Gulf of Mexico shows that corn and soybean cultivation is the largest contributor of nitrates (52%), followed by atmospheric deposition sources (16%); whereas the total phosphorus entering the Gulf of Mexico can be attributed primarily to animal manure on pastures (37%), corn and soybeans (25%), other crops (18%), and lastly urban sources (12%).

## **2.4 Statistical Load Estimation Models**

Over the past few decades, different statistical estimators have been used to predict nutrient and sediment loads in streams and rivers with sparse concentration datasets. Major load estimation approaches can be largely classified into four types: (1) ratio-based estimation; (2) average-based estimation; (3) period-weighted estimation; and (4) regression-based estimation.

### **2.4.1 Ratio-based Estimation**

Beale (1962) first introduced this class of statistical load estimation models namely the ratio based load estimation models. This approach was mandated for use in Great Lakes loading

calculations by the International Joint Commission (Richards and Holloway, 1987). The Beale ratio estimator is expressed as:

$$L_d = Q' \frac{l'}{q'} \left( \frac{1 + \frac{1}{n} \frac{S_{lq}}{l'q'}}{1 + \frac{1}{n} \frac{S_{qq}^2}{q'^2}} \right) \quad (1)$$

Where,

$L_d$  denotes estimated daily loading rate

$Q'$  denotes the mean daily flow for the year

$l_i$  and  $q_i$  denote the daily load for the day on which concentrations was determined and individual measured flows respectively

$l'$  and  $q'$  denote the mean daily loads and flows for the days on which concentrations were determined respectively

$n$  denotes the number of days on which concentrations are determined

$S_{lq}$  denotes the covariance of flux and flow

$S_{qq}$  denotes the variance of flow based on days which concentration was measured

In this technique, daily loads are calculated on days when samples are collected for evaluating nutrient concentrations (sampled days), by multiplying nutrient concentrations with the corresponding daily flow rates, and the mean daily load is then calculated for the complete monitoring period. The loads for days when no samples are collected for evaluating nutrient concentrations (unsampled days) are obtained by multiplying this mean daily load by the ratio of the mean daily flow rate for the complete monitoring period and the mean daily flow rate for sampled days. If much variation in flow rates occurs over the monitoring period, accuracies can be improved by stratifying the complete dataset based on discharge. Stratum loads can be computed by the same methodology and summed up to determine the total load for the whole monitoring duration. Ratio estimators actually use the mean of daily flows to estimate loads, and thus, they are more suited for monitoring programs with abundance of flow information for a stream but relatively scarce concentration information (Dolan et al., 1981).

Dolan et al. (1981) applied the Beale ratio estimator (BRE) to estimate total phosphorus loading in the Grand River, Michigan for a watershed draining an area of approximately 13,550 km<sup>2</sup>. They calculated total phosphorus loading for a period of one year from March, 1976-February, 1977. They estimated annual total phosphorus loads using randomly generated subsets of 25 samples. The BRE was the most accurate amongst various unbiased estimators they tested, using the root mean square error (RMSE) as a measure of goodness of fit.

Richards and Holloway (1987) carried out a similar study to estimate tributary loads with three sampling patterns, four sampling frequencies and two calculation procedures. They estimated tributary loads for three Ohio tributaries namely Maumee River, Sandusky River and Honey Creek having drainage areas ranging from 386-16,699 km<sup>2</sup>. Their results indicated that the BRE load estimates for both flow-stratified and unstratified scenarios were the most precise ones.

Burn (1990) reported that a simplified version of the ratio estimator introduced by Cochran (1977) performed indistinguishably as compared to the BRE. Ratio estimators assume the existence of a positive relationship, passing through the origin, between daily load and daily flow values. These estimators are usually considered to be the most precise unbiased estimators, and a linear relationship is typically assumed, if the variance of the daily loads is proportional to the magnitude of the daily flow rates. Ratio-based nutrient estimators have been used by a number of researchers over the past several decades (Beale, 1962; Richards and Holloway, 1987; Coats et al., 2002; Preston, 1989; Guo et al., 2002).

#### **2.4.2 Average-based Estimation**

In this simple technique, parameters on unsampled days are assigned the average value of the corresponding parameter on sampled days. If sampling frequency is variable over time, weighting has to be performed to get more accurate annual estimates. Averaging estimates can be computed for flow, nutrient concentrations, or loads.

Generally for load estimation purposes, averaging is used to compute representative concentrations, and consequently, loads are computed by multiplying these averaged concentrations with the corresponding flow rates (Walker, 1996; Short, 1999).

Dolan et al. (1981), Ferguson (1987) and Verhoff et al. (1980) developed and tested a wide array of averaging estimators. Preston et al. (1989) summarized various different average based load estimators that were developed earlier and tested them to estimate a range of tributary loads. In general, the averaging techniques can be classified under three broad categories, based on the parameter averaged:

- (1) Daily flow data with average monthly/quarterly concentrations
- (2) Average flow data
- (3) Average load data

Preston et al. (1989) studied the performance of various different averaging estimators for total phosphorus in the Grand River (1976/1977), and total lead (1977), total zinc (1978), and the PCB Aroclor-1242 (1979) in the Saginaw River. Preston et al. (1989) reported that amongst the averaging estimators that use daily flow with monthly or quarterly average concentrations had the least mean square errors (MSE). They also reported that averaging estimators that used average flow data yielded higher MSE comprising of high bias, high variance or both. The average load estimators also performed poorly and yielded higher MSE, although their precision could be improved by stratification.

### **2.4.3 Period-weighted Estimation**

This simple technique has been mainly used to estimate missing data points in a time series. For the load problem, measured nutrient concentrations of collected samples are used to estimate loads for the complete monitoring duration. There are two basic approaches to estimate missing data using period-weighting estimation:

- (1) Piecewise linear interpolation of the measured data points (Shih et al., 1998)
- (2) Use of a step function assuming that the measured concentration is the midpoint of a time-interval and the concentration is constant for the whole interval

Loads can be calculated for the entire monitoring period by multiplying the estimated concentrations with their corresponding flow rates and summing these values. Accuracy of load

estimation using period-weighted techniques is directly proportional to the number of measured concentration points.

#### **2.4.4 Regression-based Estimation**

Over the years, regression approaches have been consistently used to estimate missing nutrient and sediment concentrations in streams and rivers. Walling (1977) stated that regression approaches have traditionally been applied for estimating loads of suspended solids and other constituents. A simple- or multiple-regression relationship is developed between the dependent variable, i.e., concentration, and one or more independent variables such as flow and time. Usually as flow and concentration time series are assumed to follow a bivariate log-Gaussian (lognormal) distribution, the relationship is developed in log-log space. The back transformation of the concentrations to regular space induces a bias (Ferguson, 1986), necessitating the use of various bias correction techniques (Duan, 1983; Cohn et al., 1989).

Cohn et al. (1989) analyzed the issue of retransformation bias induced by converting the regression estimates from the “log space” to the “real space” (Ferguson, 1986). The biased regression estimator was used by Ferguson (1986), Koch and Smillie (1986), Richards and Holloway (1987) and Young et al. (1988) for estimating loads for a single stream. According to Cohn et al. (1989), in such scenarios the use of a biased regression estimator can be justified due to its simplistic approach, but in scenarios having multiple tributaries, although the random errors in the estimates for each tributary will tend to offset one another, the bias errors tends to accumulate. They developed an unbiased estimator called Minimum Variance Unbiased Estimator (MVUE), which consistently yielded comparable or better load estimates than the traditional unbiased regression estimator for a wide spectrum of conditions, especially when the sample sizes were small or loads were estimated during high flow scenarios.

Cohn et al. (1992) validated the bias correction obtained using MVUE, in load estimates from some of the major tributaries to Chesapeake Bay. According to their results the biased regression estimator gave reasonably accurate estimation of nutrient loads but a statistical significant albeit not substantial lack of fit was observed. The MVUE load estimator was found to yield satisfactory estimates of nutrient loads with statistically insignificant lack of fit.



Duan (1983) proposed a nonparametric bias correction factor and the resulting estimator is generally referred to as the smearing estimator. The smearing estimator can be expressed as:

$$L_{SM} = L_{RC} \frac{1}{M} \sum_{i=1}^M \exp[e(i)] \quad (2)$$

Where,

$L_{SM}$  denotes load estimates from the smearing estimator

$L_{RC}$  denotes load estimates from the rating curve (regression) estimator

$e(i)$  denotes the regression residual

$\exp$  denotes the exponential function

Guo et al. (2002) comprehensively analyzed the performance of three major classes of load estimators, the biased rating curve (regression) estimator, ratio estimator and the flow-weighted average estimator. Additionally they also assessed the performance of the MVUE and smearing estimator bias correction techniques. For different combinations of sampling frequencies and monitoring durations, Guo et al. (2002) compared “true” nitrate-N load estimates from the Upper Sangamon River watershed in Illinois with estimates from the above mentioned estimators and bias correction techniques. Their results indicated that bias correction wasn’t required in this scenario and both the bias correction techniques actually increased the bias. More stream nutrient concentration datasets from a wide spectrum of locations are needed to quantify the performance of various statistical load estimation approaches. There are only a handful of studies which have been conducted over the past few decades that compare the performances of various nutrient load estimation models. Barring a few, most of these studies are hampered by a very limited duration of continuous concentration datasets which is used to compare the performances of the different load estimation models. Also, many of the estimators and bias correction techniques vary in performance for datasets with different characteristics and from different locations.

## 2.5 Composite Method

Aulenbach and Hooper (2006) proposed an alternative method, called the composite method to estimate solute loads. They used an extensive dataset collected at the outlet of the Panola Mountain Research Watershed (PMRW) near Atlanta, Georgia, USA, to illustrate this

method, which combines the strengths of regression and period-weighted approaches. Residual concentrations are computed for all sampled days by subtracting concentrations estimated using the regression model from the corresponding measured values, and a piecewise continuous linear function of the residual concentrations over time is developed. This residual concentration from this composite function is then subtracted from the regression concentrations on unsampled days.

Aulenbach and Hooper (2006) demonstrated that the composite method improved load estimation accuracies over short time intervals and allowed for better trend analysis of load estimates. The composite method was based on the principle of autocorrelation amongst the residual concentrations, and it used a piecewise linear interpolation to distribute residuals in between sampled concentrations. This concept is an advanced form of curve fitting which has created a new class of error correction techniques for load estimation models. Based on this concept, by optimizing the distributions used to assign residuals in between sampled concentrations potentially more accurate new error correction techniques may be developed.

# **CHAPTER 3**

## **OBJECTIVES**

The overall goal of this study was to improve nitrate-N load estimation accuracies for scenarios with a wide spectrum of monitoring durations and sampling frequencies using various load estimation models and error correction techniques. Specifically, the objectives of this study were to:

- 1) Confirm and validate the findings of previous research related to changes in load estimation accuracies by increasing monitoring duration and sampling more frequently.
- 2) Evaluate the performance of various load estimation models for different monitoring scenarios.
- 3) Develop and analyze the performance of various error correction techniques for load estimation models, and assess if they can be used to design cost-effective monitoring plans, requiring fewer observations for the same load estimation accuracies.
- 4) Test if incorporation of monitoring thresholds in conjunction with fixed monitoring can be used to design cost-effective monitoring plans, having the lowest load estimation errors for any fixed number of samples collected.

# **CHAPTER 4**

## **SITE AND DATASET DESCRIPTIONS**

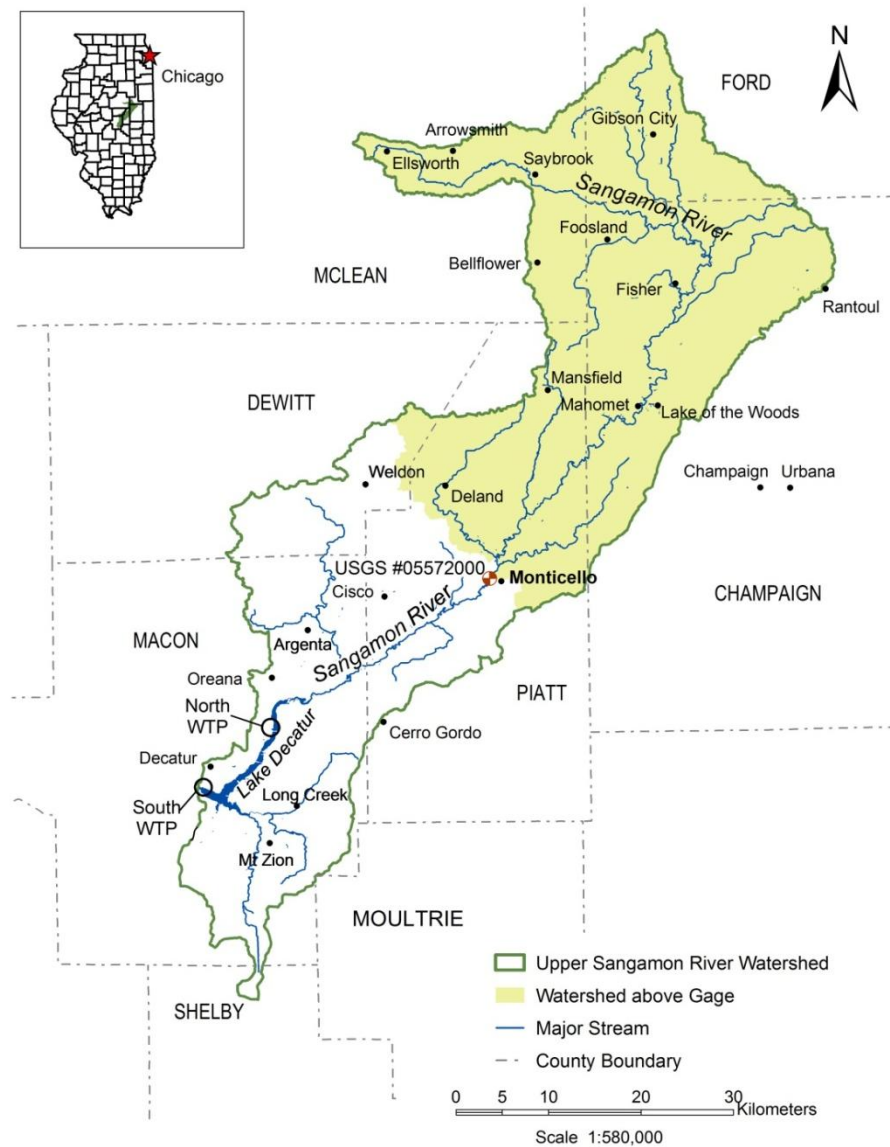
### **4.1 Watershed Descriptions**

Two watersheds in north-central Illinois, the Upper Sangamon River at Monticello (United States Geological Survey [USGS] 8-digit HUC number 05572000), and Vermilion River at Pontiac (USGS 8-digit HUC number 05554500), were monitored over a period of six years from May 1993 to April 1999. The watersheds have similar drainage areas, and the topography of both watersheds is characterized as being almost flat to gently sloping, tilled farmland. Soil types of both watersheds are dominated by soils with poor natural drainage, which has led to the extensive adoption of tile drain systems in the watersheds. Agriculture is the primary land use for both watersheds, with corn-soybean rotation being the most common cropping pattern. Annual precipitation trends for both watersheds indicate that precipitation ranges from 890 to 1,015 mm annually, most of which occurs in spring and summer. As seen in most of the watersheds in Illinois, both watersheds in this study display well-defined seasonal flow variability. Flows during the spring (March–June) are expected to be the highest with lowest flows generally occurring through late summer and fall. Most major water quality issues for both watersheds used in this study relate to large-scale sedimentation and nonpoint nutrient pollution, mainly nitrate-N and phosphorus (IEPA, 2007; IEPA, 2009).

### **4.2 Dataset Description**

#### **4.2.1 Upper Sangamon River at Monticello**

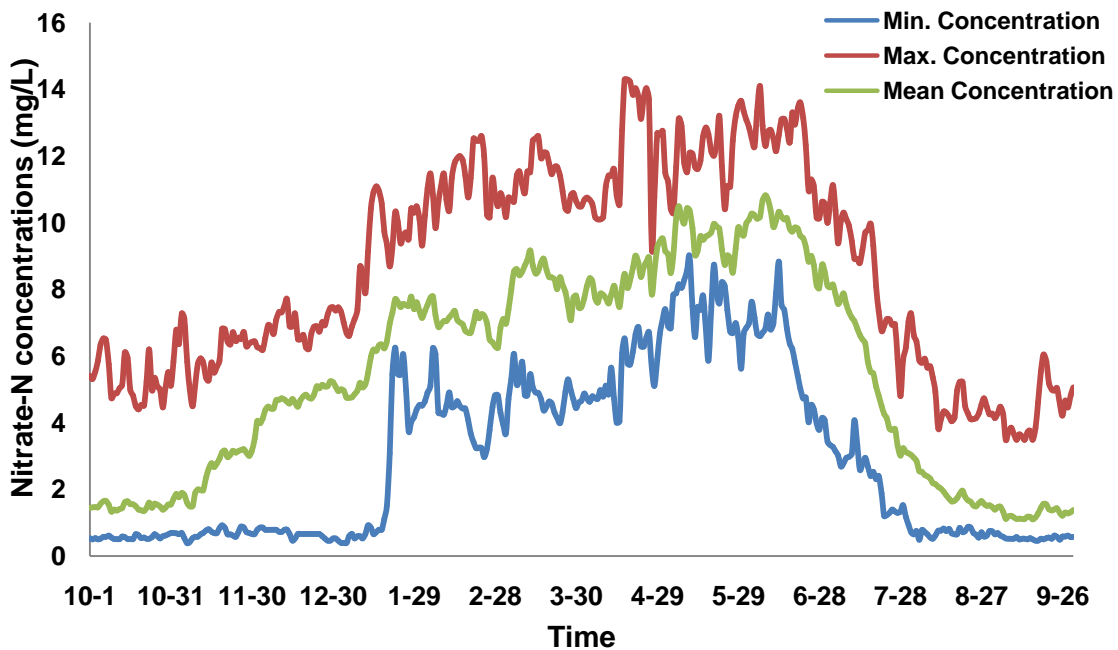
The Upper Sangamon River watershed is located in east-central Illinois and drains an area of approximately 3161 km<sup>2</sup> (Figure 1). This study specifically focuses on a sub-watershed of the Upper Sangamon River located upstream of the city of Monticello, draining an area of 1406 km<sup>2</sup>.



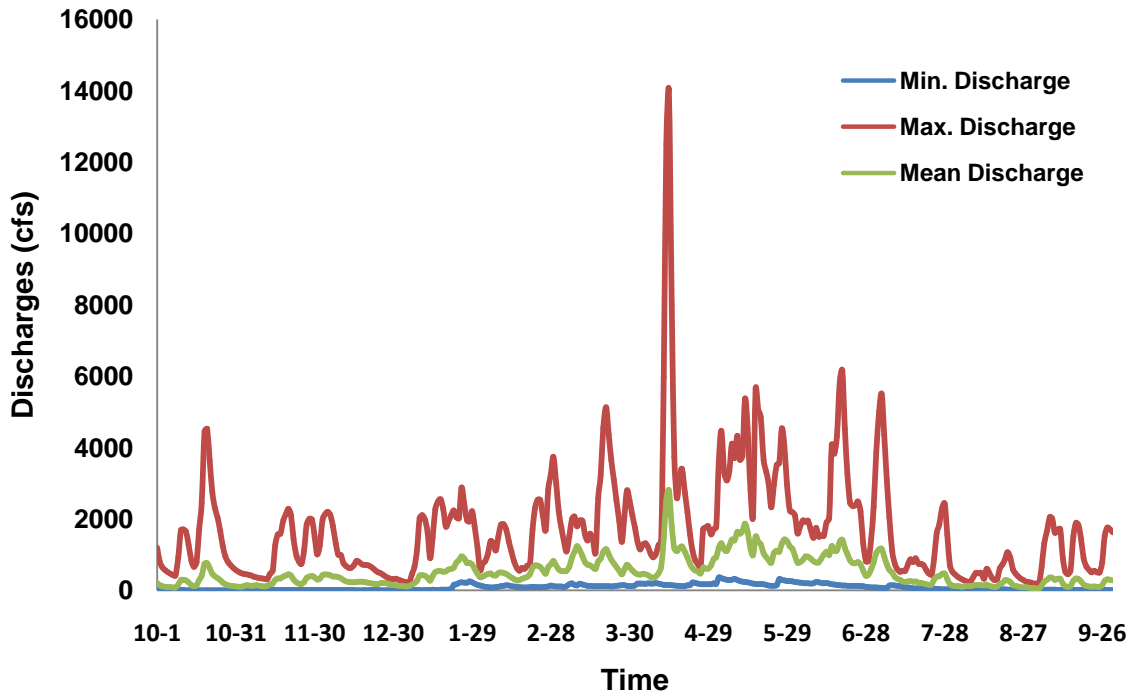
**Figure 1. Upper Sangamon River watershed.**

This watershed eventually discharges into Lake Decatur, which is located approximately 39 kilometers (km) downstream of Monticello. Over the years, the U.S. Geological Survey has operated seven continuous stream gaging stations in the Upper Sangamon River watershed. At one of these gaging stations on the Sangamon River at Monticello, they monitored discharge from the sub-watershed studied in this project. Mean daily stream flow data from this gaging station, from May 1993 through April 1999, were retrieved from the USGS and used in this

study. Since May 1993, the Illinois State Water Survey (ISWS) independently collected regular samples from this site for nitrate-N analysis (Keefer and Demissie, 2000). The city of Decatur collected instantaneous samples for nitrate-N analysis, from the north and south water treatment plant intake, daily. It was assumed that concentration values of these instantaneous samples represented mean daily concentrations. The daily nitrate-N concentration data from the north water treatment plant intake were retrieved from the city of Decatur for May 1993 through April 1999. Guo et al. (2002) used fitted stepwise regression to generate daily nitrate-N concentration data at Monticello using the daily nitrate-N concentrations at the north water treatment plant at Lake Decatur and the data collected by the ISWS. The product of the daily measured or generated nitrate-N concentrations at Monticello and the recorded mean daily flow rates yielded daily load values of nitrate-N at the Monticello gaging station. Annual load estimates for nitrate-N were obtained by summing these daily load values over a year (May 1 –April 30). The total load for the period of study was obtained by summing the annual values.



**Figure 2. Minimum, maximum, and mean nitrate-N concentrations at Monticello (Sangamon River) for 1993-1999 observation years averaged over a water year.**

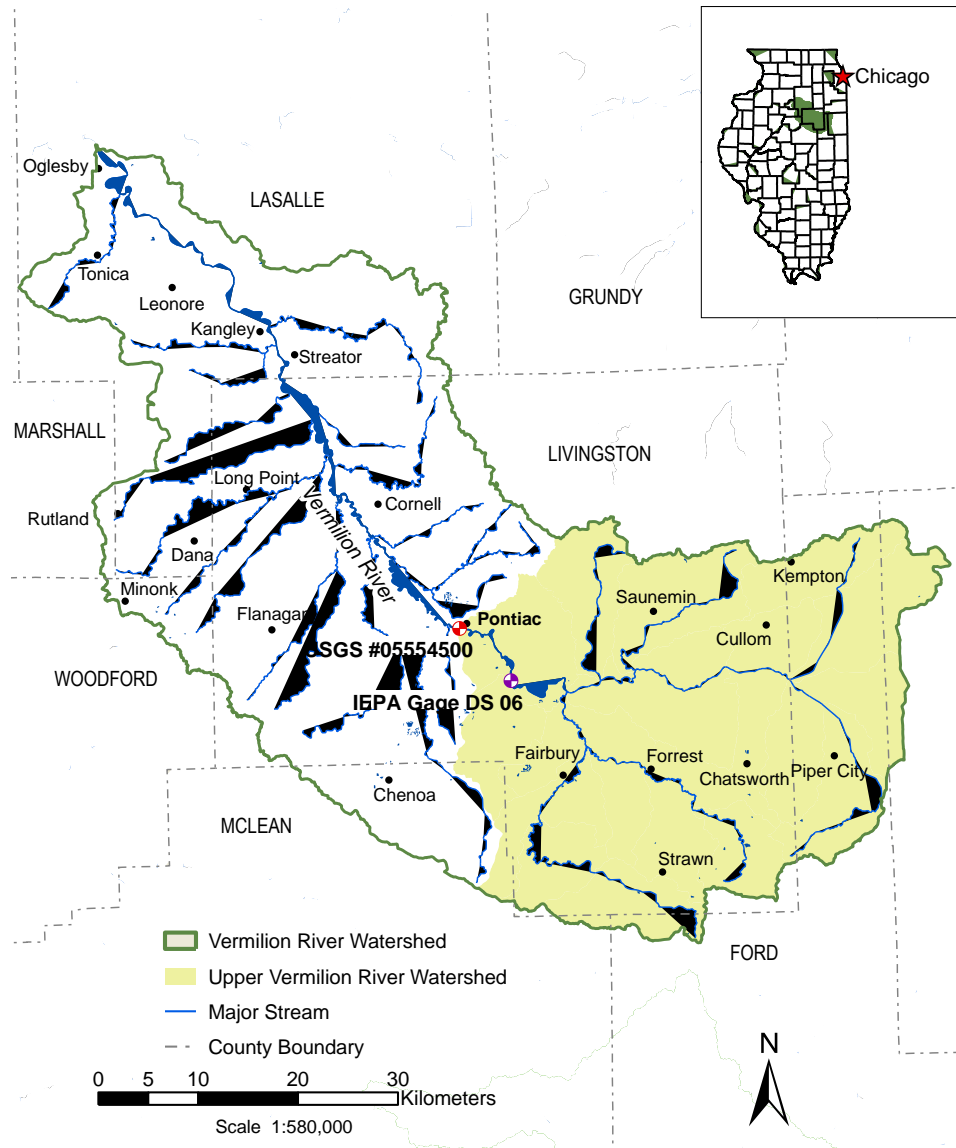


**Figure 3. Minimum, maximum, and mean discharges at Monticello (Sangamon River) for 1993-1999 observation years averaged over a water year.**

A nitrate-N yield of approximately 26.4 kilograms per hectare per year (kg/ha/year) for the six-year study period was obtained by dividing the total load by the contributing watershed area. Data from May 1993 to April 1999 were used to calculate the mean, minimum, and maximum nitrate-N concentrations and flow rates for each complete water year. Nitrate-N concentrations and flow rates for the Upper Sangamon River watershed at the Monticello gauging station are summarized in Figures 2 and 3, respectively.

#### **4.2.2 Vermilion River at Pontiac**

The Vermilion River watershed, located in north-central Illinois, drains an area of approximately 3424 km<sup>2</sup> (Figure 4). This study focuses specifically on the sub-watershed of the watershed located upstream of the gaging station at Pontiac, that drains an area of 1492 km<sup>2</sup>.



**Figure 4. Vermilion River (Illinois basin) watershed.**

Mean daily discharge data at the Pontiac gaging station, collected by the Northern Illinois Water Corporation (NIWC) from May 1988 to April 1999, were retrieved and used in this study. To analyze nitrate-N, NIWC also collected daily samples close to a water treatment plant in Pontiac. Prior to November 1995, the cadmium reduction method was used to determine nitrate-N concentrations; subsequently the NIWC switched to the ion selective electrode method to



determine concentrations. To measure nitrate-N concentrations, the ion selective electrode method uses a stable and robust probe, which generally provides a slightly more accurate estimate of nitrate-N concentrations as compared to the cadmium reduction method (Capelo et al., 2007). A major concern regarding the quality of the concentration data measured was that the samples were collected from the edge of the river upstream from the treatment plant, with no attempt being made to collect multiple samples across the depth and width to obtain more representative estimates of nitrate-N concentration. However, the Illinois Environmental Protection Agency (IEPA) operated a sampling site (McDowell) 14.5 km upstream from the Pontiac gaging station, collecting nine comprehensive depth- and width-integrated samples annually, to estimate representative nitrate-N concentrations across the Vermilion River. For the complete monitoring period from May 1988 through April 1999, a high degree of correlation was observed between the nitrate-N concentration measurements for samples collected by NIWC and IEPA, with the coefficient of determination ( $R^2$  value) being 0.79 (Figures 5 and 6).

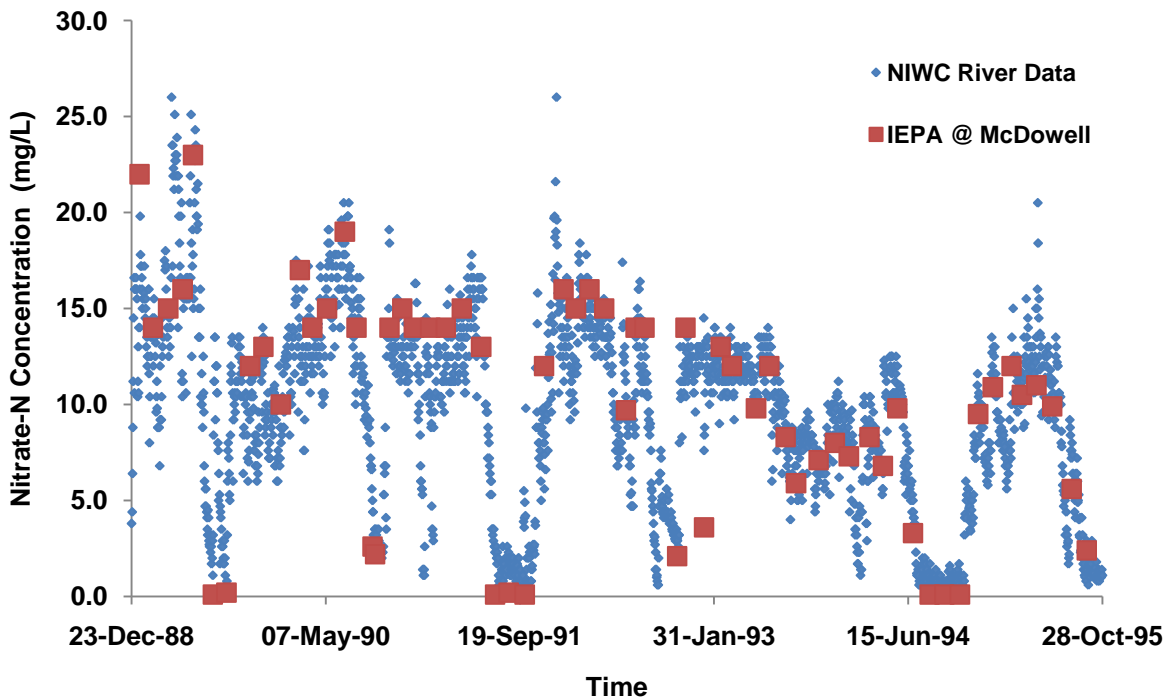
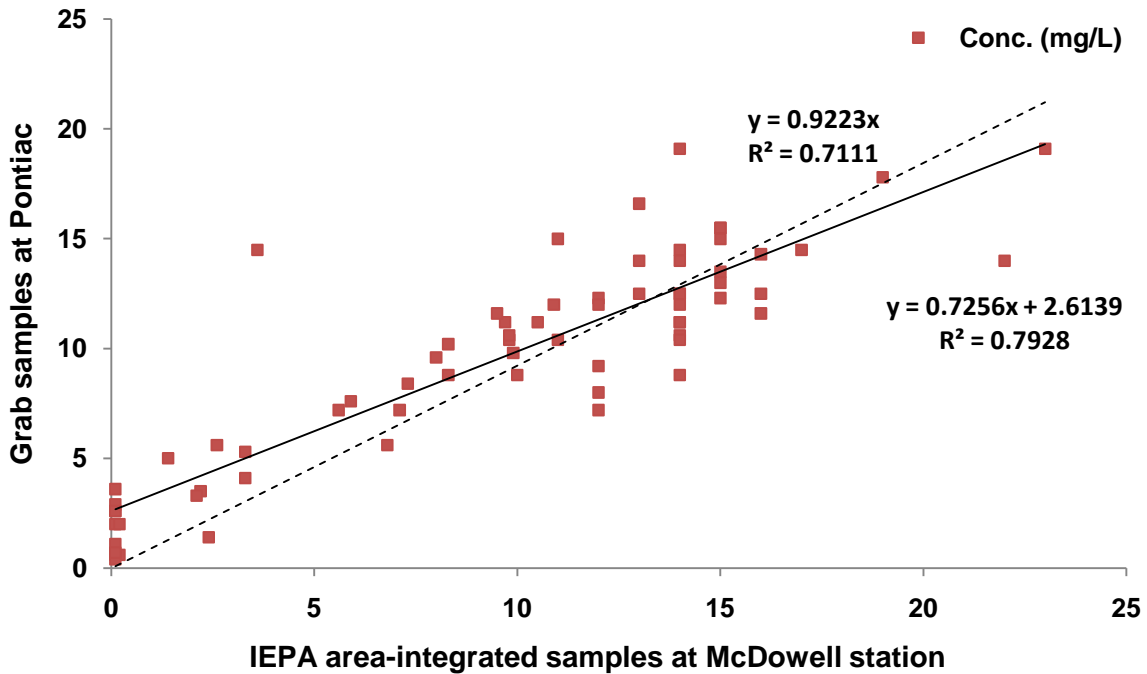


Figure 5. Nitrate-N concentration data at McDowell (area-integrated) collected by IEPA and Pontiac (grab from shore) collected by NIWC.



**Figure 6. Correlation of nitrate-N concentration data at McDowell (area-integrated) collected by IEPA and Pontiac (grab from shore) collected by NIWC.**

Thus, the nitrate-N concentration measurements taken by NIWC at Pontiac were assumed to be accurate for load estimation throughout the complete monitoring duration. Daily nitrate-N load estimates at the Pontiac gaging station were calculated by multiplying daily measured nitrate-N concentrations and recorded mean daily discharges. Annual nitrate-N load estimates were obtained by summing daily load values over a year. To maintain consistency with the Upper Sangamon River at Monticello, only the nitrate-N load estimates from May 1993 through April 1999 were used in this study.

The average nitrate-N yield for the watershed for the six-year period of study was approximately 23.4 kilograms per acre per year (kg/acre/year). Nitrate-N concentrations and flow rates, respectively, for the Vermilion River watershed upstream of the Pontiac gaging station are summarized in Figures 7 and 8.

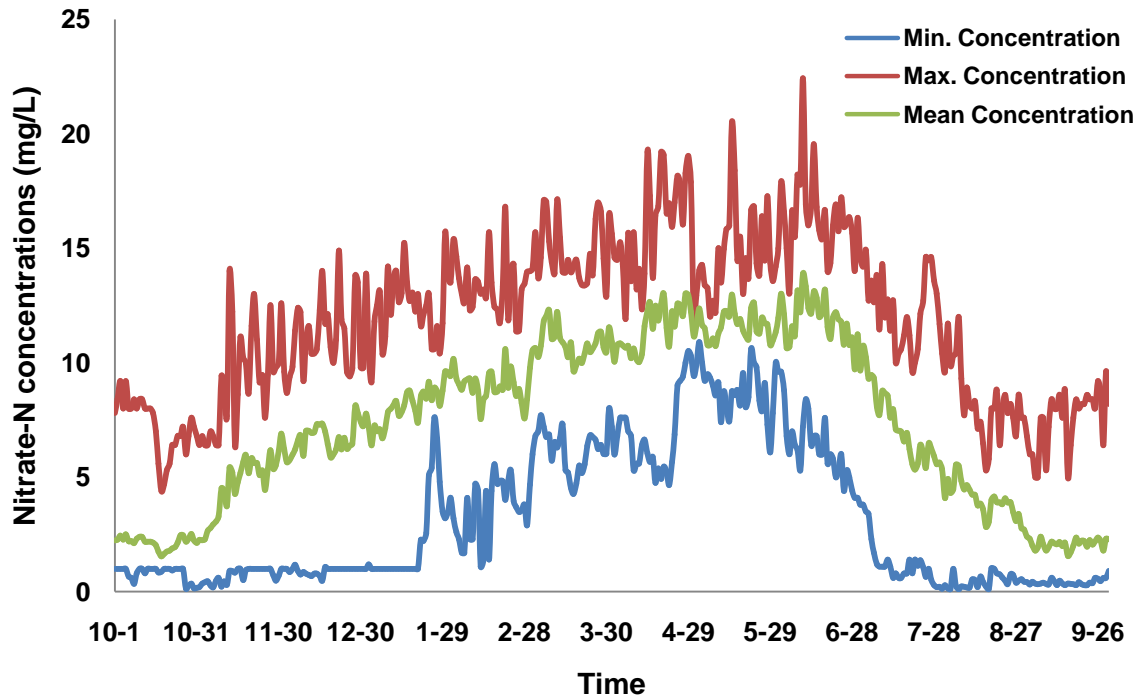


Figure 7. Minimum, maximum, and mean nitrate-N concentration at Pontiac (Vermilion River) for 1993-1999 observation years averaged over a water year.

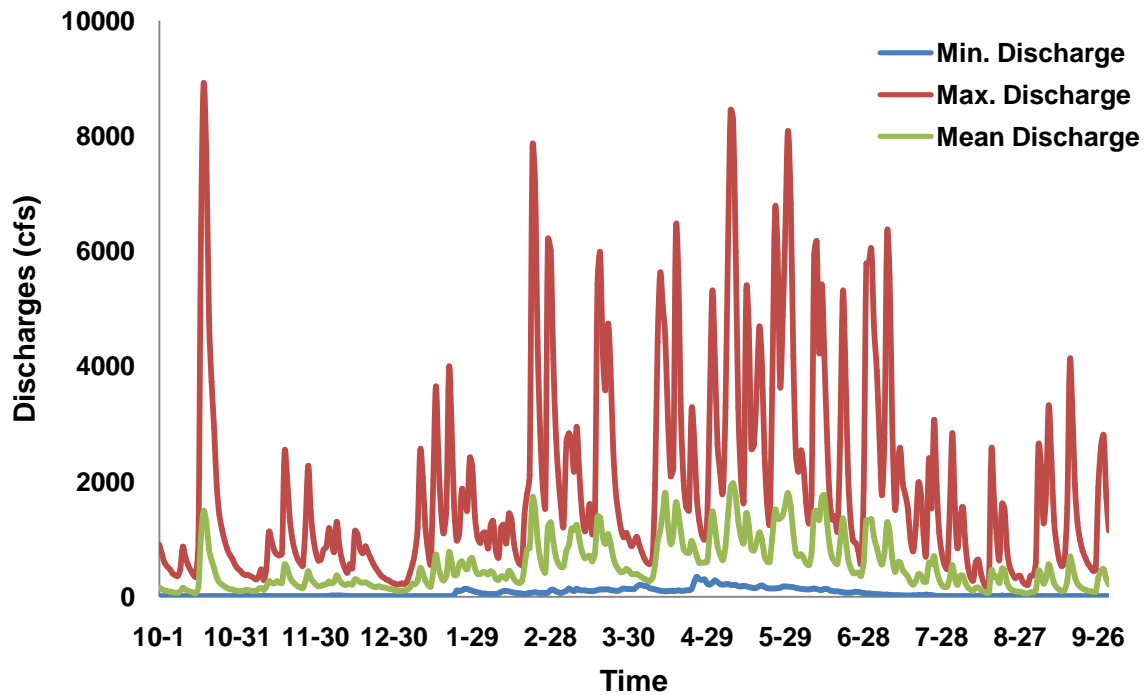


Figure 8. Minimum, maximum, and mean discharges at Pontiac (Vermilion River) for 1993-1999 observation years averaged over a water year.

### 4.3 Flashiness Index

Baker et al. (2004) developed a criterion called the flashiness index to characterize flow regimes of various Midwestern streams. Flashiness reflects frequency and rapidity of short-term changes occurring in stream flow regimes. A stream hydrograph with higher storm peaks and lower base flow values is flashier than a stream with higher and more stable base flows and smaller storm peaks. The R-B Index was computed individually all six years of study for both watersheds. The R-B flashiness index values for flows from the Vermilion River sub-watershed at Pontiac were consistently significantly higher than that for flows from the Upper Sangamon sub-watershed at Monticello. A graph of different flashiness index values for both watersheds used in this study over the complete monitoring duration (1993-1999) is shown in Figure 9.

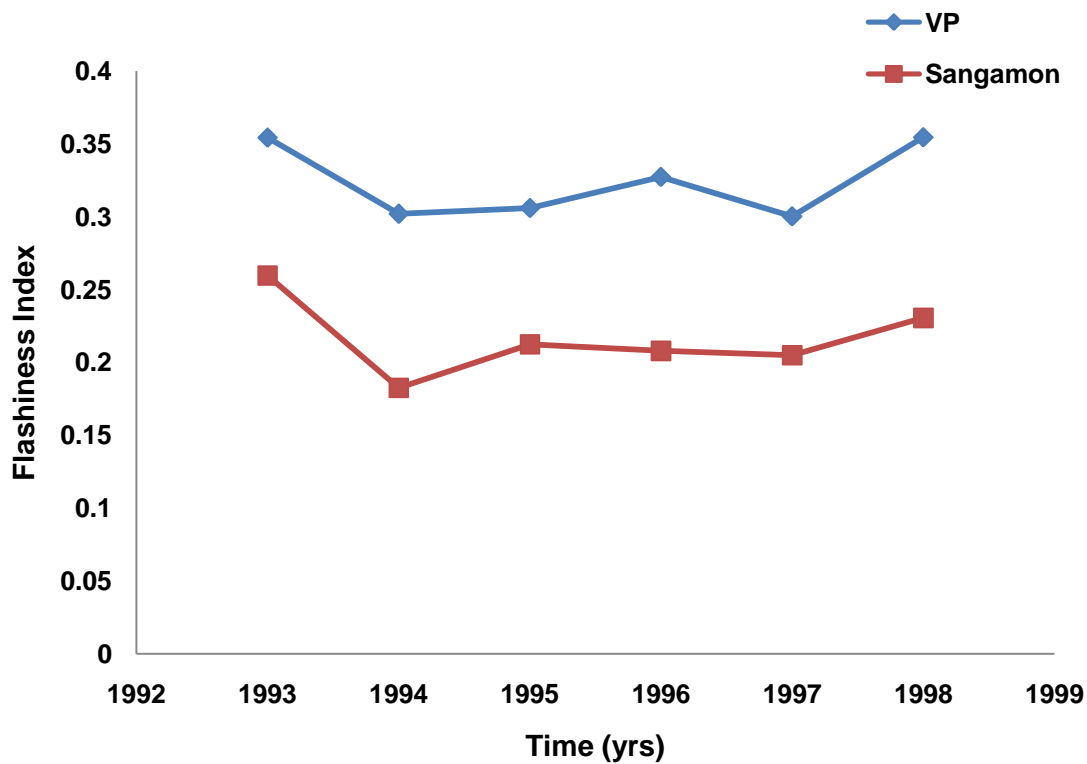
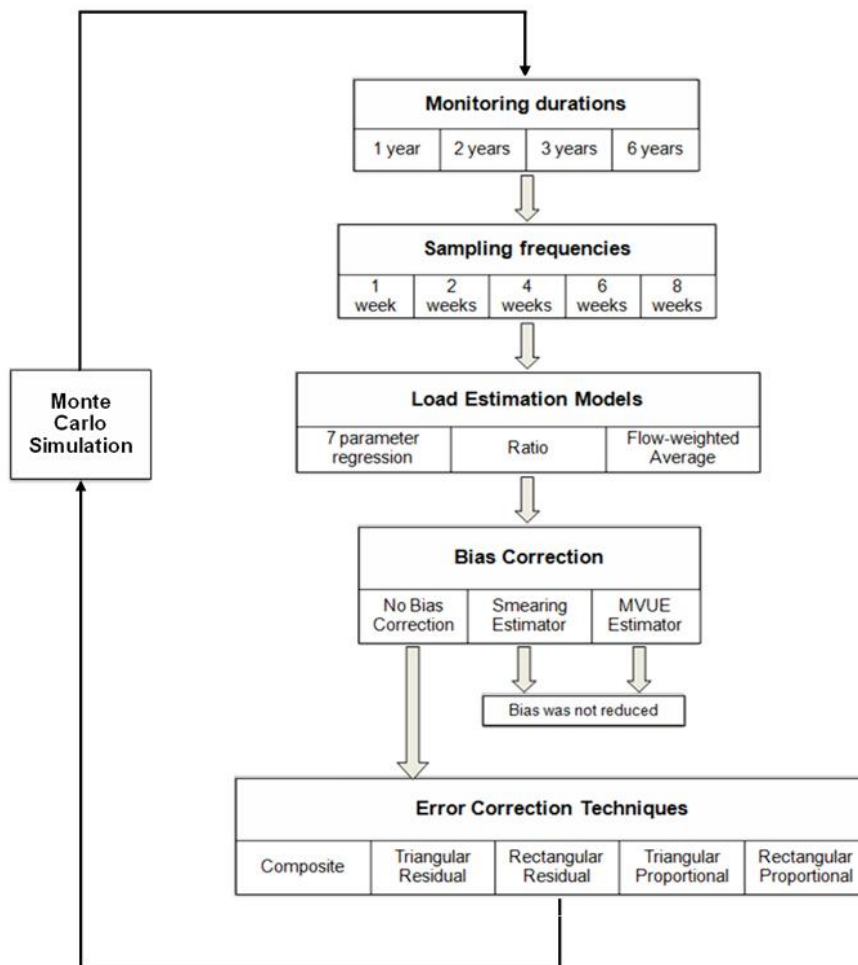


Figure 9. Richards-Baker flashiness index for Vermilion and Sangamon Rivers for 1993-1999 observation years.

# CHAPTER 5

## METHODOLOGY

The approach used in this study to analyze the performance of different load estimation models and error correction techniques is shown schematic in Figure 10. Unique combinations with different monitoring durations (1 year-6 years), sampling frequencies (1 week-8 weeks), and three different load estimation models were tested. The composite method along with four new error correction techniques were applied to each of the unique scenarios. 700 Monte Carlo iterations were performed to compare load estimates to the “true” loads from both the watersheds used in the study.



**Figure 10. Schematic of project approach.**

## 5.1 Load Estimation Methods

Load estimation is commonly a vital objective in most of the water quality monitoring programs. Mass fluxes or loads of a nutrient over a certain monitoring period across any channel can be calculated using continuous flow rate and concentration measurements. The total solute load is can be estimated by integrating solute concentrations and discharges over time. Typically in monitoring programs discharge is measured often and the concentration is measured less frequently. To have accurate load estimates there is a need to estimate concentrations between relatively infrequent concentration samples. Various different approaches have been studied and used to estimate missing concentrations. Most monitoring programs are designed to measure concentrations at weekly to eight weekly intervals, and to use statistical models to expand these values into a continuous (daily) concentration dataset. In this study, the following three statistical load estimation models were compared to estimate missing nitrate-N concentrations:

### 5.1.1 USGS Seven-Parameter Estimator

The USGS seven-parameter estimator developed by Cohn et al. (1992), also known as the rating curve method, is the most widely used regression-based load estimator. It is a multiple non-linear regression model which develops and utilizes a relationship between a sparsely collected dependent variable (i.e., concentration) and continuous independent variables such as flow rate and time of year. To achieve a better linear relationship between the dependent and independent variables, and to reduce the influence of extreme flows on load calculations, concentrations are log transformed in the model. The model can be written in the following form:

$$\ln(C) = \beta_0 + \beta_1 \ln\left(\frac{Q}{Q'}\right) + \beta_2 \left[\ln\left(\frac{Q}{Q'}\right)\right]^2 + \beta_3(T - T') + \beta_4(T - T')^2 + \beta_5 \sin(2\pi T) + \beta_5 \cos(2\pi T) + \varepsilon \quad (3)$$

Where,

$\ln(\ )$  denotes natural logarithmic function

C and Q denote the nitrate concentration and flow rate

T denotes time measured in years

$\varepsilon$  denotes the errors, which are assumed to be normally distributed with a mean of 0 and a variance of  $\sigma_\varepsilon^2$

$\beta_0, \beta_1 \dots \beta_6$  denote the various regression parameters estimated from the available dataset

Q' and T' denote centering variables used in the model.

Concentrations from sampled days, along with the corresponding mean daily flow rate and Julian days are used to compute the seven parameters. The model is then used to estimate log-concentrations, and ultimately concentrations, on unsampled days, and the expanded concentration data set is then used to obtain load estimates for any subset of the monitoring period:

$$L_{RC} = \sum_{j=1}^N C_j Q_j \Delta T \quad (4)$$

Under the assumption that regression residuals are independent and normally distributed, Koch and Smillie (1986) and Ferguson (1987) found that the regression model under predicts concentrations (i.e., a bias is observed in the concentrations). This bias can be attributed to application of a regression model to log-transformed data. To account for the bias in concentration predictions, Cohn et al. (1992) proposed the application of a minimum variance estimator (MVUE). Accurate load estimations of nutrient loads discharging into the Chesapeake Bay were reported by (Cohn et al., 1992).

$$L_{MVUE} = L_{RC} g_m \left[ \frac{m+1}{2m} (1 - v) s^2 \right] \quad (5)$$

Where,

$m$  denotes the number of observations used to calibrate the model minus the number of parameters in the model

$v$  denotes a leverage term, which is a function of independent variables for which the load concentration is calculated

$s^2$  denotes the estimated variance of the regression residuals

$g_m$  is function defined as:

$$g_m = \sum_{p=0}^{\infty} \frac{m^p(m+2p)}{m(m+2)\dots(m+2p)} \left(\frac{m}{m+1}\right) \left(\frac{z^p}{p!}\right) \quad (6)$$

Wang and Linker (2008) demonstrated that adding additional parameters to the original seven-parameter model can better account for hysteresis in sediment transport in a regression model and thereby improve load estimation accuracies for scenarios with frequent sampling and/or seasonal variations.

### 5.1.2 Ratio Estimator

Cochran (1977) developed the ratio estimator, a simple statistical model that is used to predict missing concentrations. Daily loads on sampled days are computed by multiplying mean daily flow rates and measured nutrient concentrations. To compute load estimates on unsampled days, the mean daily load for the sampled days is then multiplied by the ratio of the cumulative daily mean flow rates for the entire monitoring period and the daily mean flow rate for sampled days. Ratio estimators are, therefore, primarily based on the assumption that the ratio of load to flow rate remains constant for sampled days and unsampled days through the complete monitoring duration. It can be represented by the following equation:

$$L_{\text{ratio}} = \frac{l'}{q'} \sum_{j=1}^N Q_j \quad (7)$$

Where,

$l'$  and  $q'$  denote the means of daily loads and flow rates, respectively, for sampled days

To increase precision, the dataset can be stratified on the basis of daily flow rates, and separate loads estimated for each stratum used the methodology outlined above. Total load estimates for the entire monitoring period can then be calculated by summing the stratified loads. This ratio estimator was used in this study and various correction methods were applied to it to gauge the improvement in load estimation accuracies, if any.



### 5.1.3 Flow-Weighted Average Estimator

Walker (1996) suggested a flow-weighted average estimator, to estimate nutrient and sediment loads from various watersheds in Illinois, that was subsequently used by Short (1999). This technique required that the complete dataset be stratified based on flow. In this study, the 1993-1999 dataset was subdivided into three strata, similar to studies done by Short (1999) and Guo et al. (2002) where datasets were divided into:

- (1) All flow rates less than half of the mean flow of the complete monitoring period ( $Q_{\text{mean}}$ );
- (2) Flow rates between half and twice the mean flow; and
- (3) Flow rates greater than twice the mean flow.

The type of flow-weighted average load estimation model used in this study is essentially a stratified version of the standard ratio-estimation model. Stratified loads were computed using the following relation:

$$L_{\text{AVE}} = \frac{\sum_{i=1}^{M_1} C_i Q_i}{\sum_{i=1}^{M_1} Q_i} \sum_{j=1}^{N_1} Q_j \quad (8)$$

Where,

$C$  denotes the nitrate concentration

$Q$  denotes the flow rate

$M_1$  denotes the number of calibration observations for the first strata

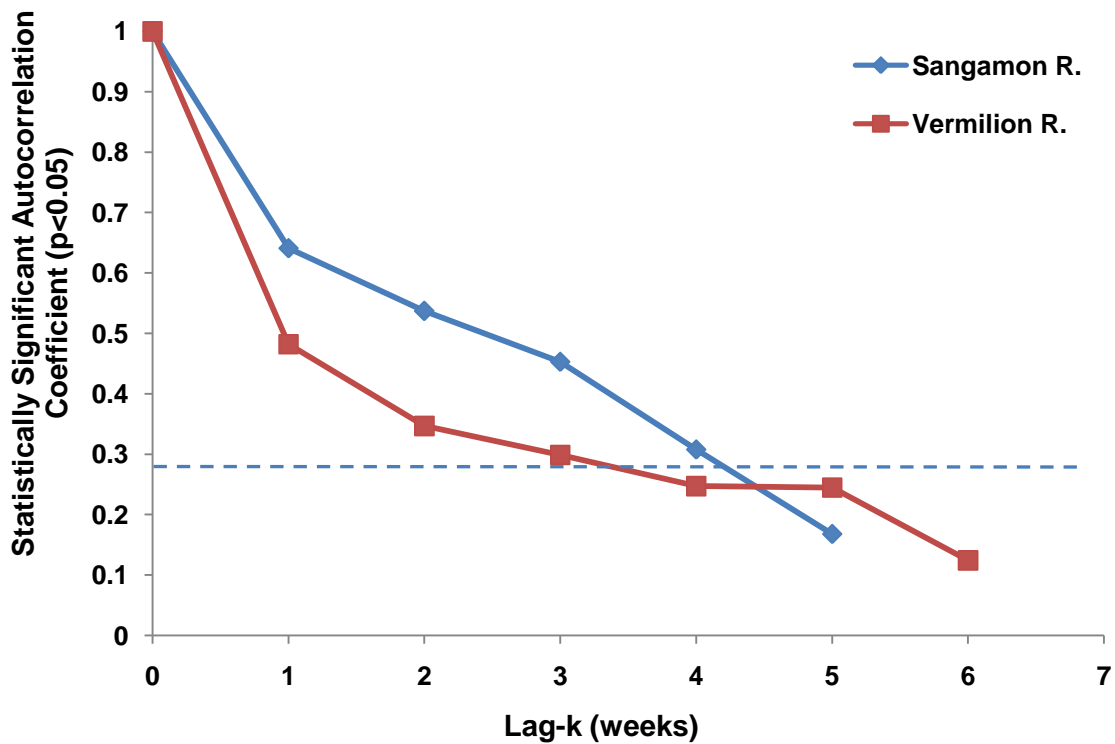
$N_1$  denotes total number of observations respectively for the first strata

Total loads for the complete monitoring duration were computed by summing up the individual stratified loads.

## 5.2 Autocorrelation of Modeling Residuals

Regression-based approaches to estimate missing nutrient concentrations generally fail to account for any autocorrelation structure that may exist in residual concentrations.

Autocorrelation analysis of residual concentrations can be performed to determine the presence of serial autocorrelation, which, if present, can then be used to improve estimates of nutrient concentrations. A recent USGS report states that increasing sampling frequency usually increases serial autocorrelation in residual concentrations (Aulenbach et al., 2007). It also states that serial autocorrelation over 0.2 in residual concentrations can be used improve estimations of nutrient load fluxes if it is taken into account in the estimation model.



**Figure 11. Correlogram of residual concentrations (difference between observed concentrations and modeled concentrations using seven-parameter regression model) for Sangamon and Vermilion Rivers for 1993-1999 observation years.**

This potential improvement led to the development of a new approach, called the composite method, to improve solute load estimations from regression methods (Aulenbach and Hooper, 2006). Inspired by this composite method, four error correction techniques were developed and tested in this study. The residual concentrations from the models fitted to both the watersheds used in this study exhibit the presence of significant and persistent serial autocorrelation, with coefficients greater than 0.2 for lags of up to five weeks (Figure 11).

### 5.3 Composite Method

As stated above, the composite method was proposed by Aulenbach and Hooper (2006). It combines the strengths of the traditional regression-based approach with the simple period-weighted approach that has frequently been used to predict missing stream water solute concentrations. Traditional regression-based load estimation approaches do not use residual concentrations (difference in regression model predicted and observed/measured concentrations), in the prediction missing concentrations. The composite method incorporates a new residual-load component into the conventional model load estimation model. Like period-weighted approaches, step-wise linear interpolation is used to estimate residual concentrations on unsampled days. The resulting residual concentration function is multiplied by flow, integrated over time, and then subtracted from the conventional regression model load estimation function. The composite method is expressed as:

$$L_{\text{composite}} = \int C_m(t)Q(t)\Delta t - \int C_\epsilon(t)Q(t)\Delta t \quad (9)$$

Where,

$C_m$  denotes solute concentration computed using the regression approach

$C_\epsilon$  denotes solute concentration computed using the period-weighted approach

$Q(t)$  denotes flow-rate at time  $t$

$\Delta t$  denotes the total monitoring duration for which loads are calculated.

In essence, the composite method adjusts the regression-model-predicted concentrations on sampled days to the observed values by subtracting the value from a function developed by step-wise linear interpolation of residual concentrations. In the composite method for the regression function, Aulenbach and Hooper (2006) used a hyperbolic function, in the place of the log-linear function, to develop a relationship between concentration and discharge.

$$C_M(t) = a_0 + \frac{a_1}{1 + \beta Q^*(t)} + a_2 m + a_3 \sin\left(\frac{2\pi D}{365}\right) + a_4 \cos\left(\frac{2\pi D}{365}\right) + a_5 \sin\left(\frac{2\pi D}{182.5}\right) + a_6 \cos\left(\frac{2\pi D}{182.5}\right) \quad (10)$$

Where,

$m$  denotes a dummy variable depending on the slope of the hydrograph

$D$  denotes the day of year (0 to 365)

$t$  denotes time

$\beta$  denotes a fitted hyperbolic regression model parameter to linearize data

$a_n$  (0 to 6) denote regression model parameters

## **5.4 Development of Error Correction Methods**

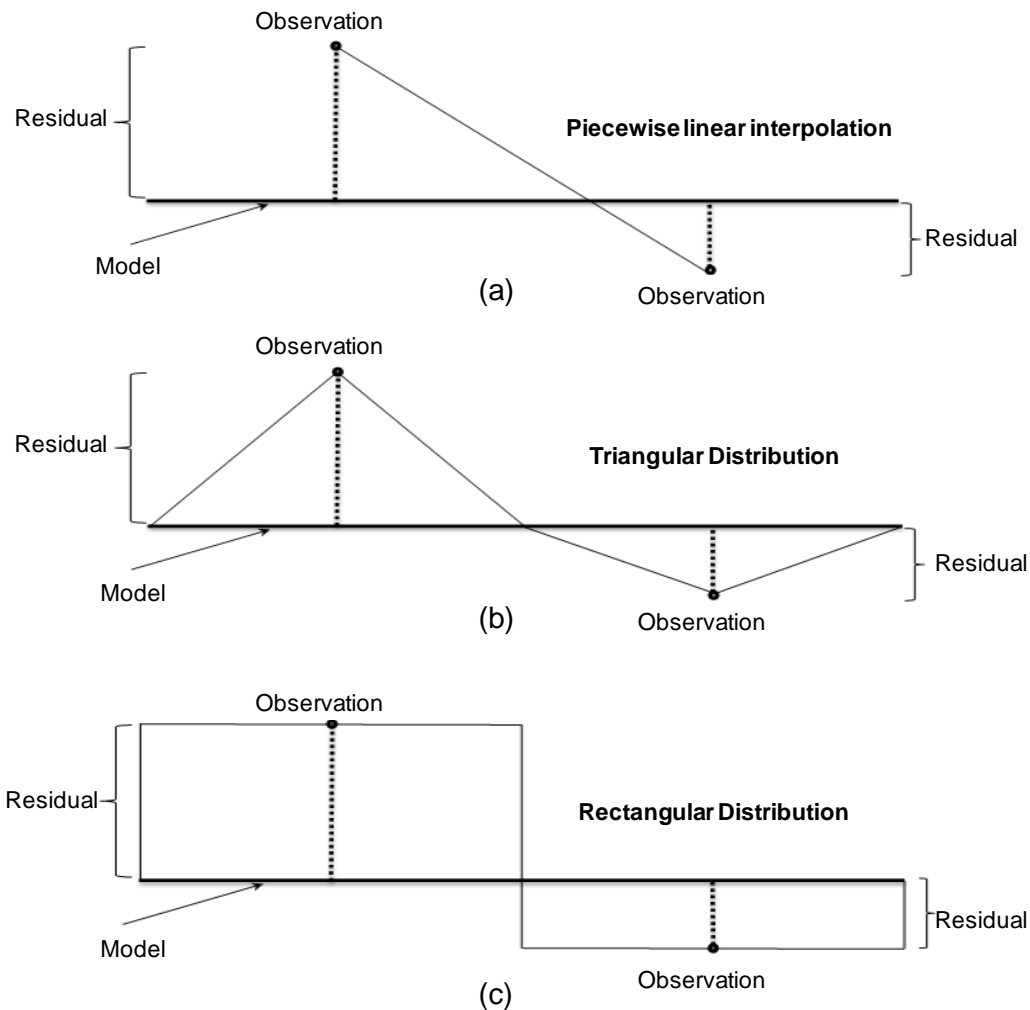
Similar to the composite method, in this study the temporal correlation in modeling errors (residual concentrations) were taken into consideration by assigning errors on unsampled days, based on known errors on proximate sampled days. Firstly, residual concentrations were calculated on sampled days, by finding the difference in regression-model-predicted concentration and observed concentrations. The mid-points of the time interval between all the pairs of consecutive sampled days were then determined. These midpoints were then set as vertices in the estimation of the residuals on adjacent unsampled days. In this study instead of using step-wise linear interpolation, as in the composite method, rectangular- and triangular-shaped distributions were used to assign errors in the vicinity of a known error. For the rectangular distribution, the residual at an unsampled day was assigned the value of the residual on the closest sampled day; while for the triangular distribution, the magnitude of the residuals were assumed to be maximum on the sampled days, and vary linearly to zero at mid-points of the time-intervals between consecutive sample days.

### **5.4.1 Concepts of Residual and Proportion**

In this study another set of correction techniques, based on proportional concentrations rather than residual concentrations, was also developed. A proportional concentration is defined as the ratio between the observed/measured concentration and the model estimated concentration on a sampled day:

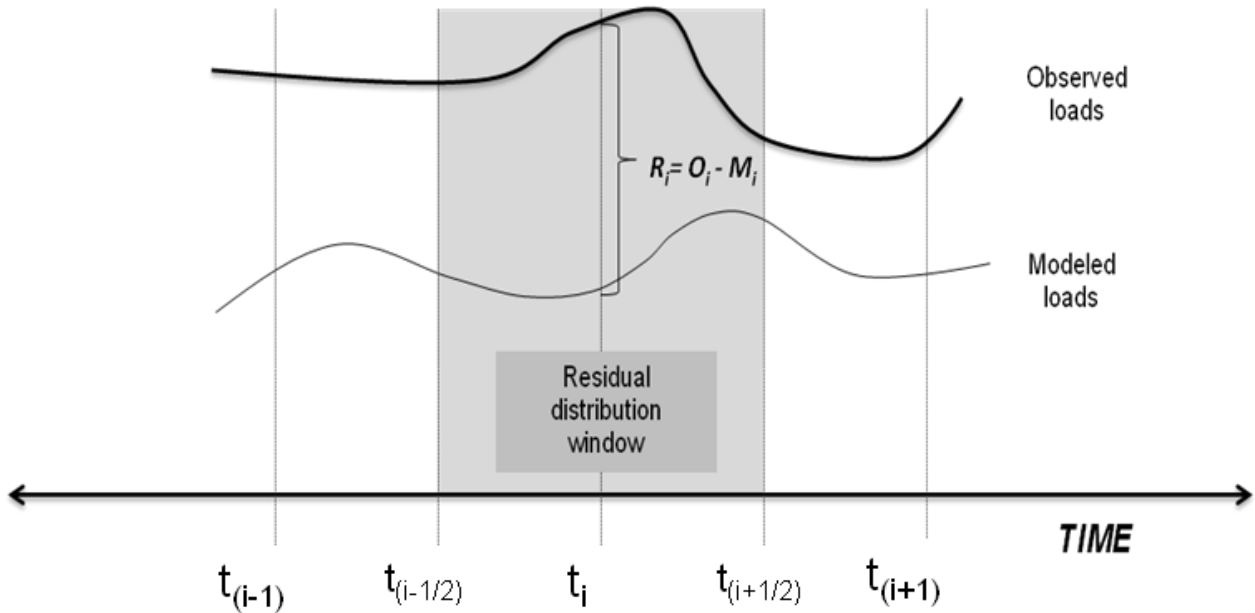
$$\text{Proportion}_N = \frac{\text{Observed concentration}}{\text{Modeled concentration}} \quad (11)$$

Once the set of proportional concentrations was obtained for all sample observations, the rectangular and triangular distributions again were used to assign concentrations in the vicinity of known proportional concentrations. Schematic diagrams of the different distributions that were used to develop continuous functions, based on residual concentrations, are shown in Figure 12.



**Figure 12. Schematic of the various distributions for residual concentrations, a) composite method (piecewise linear interpolation); b) triangular distribution; c) rectangular distribution.**

Figure 13 shows a schematic of the extent of the residual distribution window. Residuals for days in between the sampled days were calculated using the triangular and rectangular (step) distributions and also by step-wise linear interpolation of the known residual values. Known residuals were used to calculate the missing residuals for a window which extended equally on either side of the known residual up to the midpoint between two known adjacent residual values. In Figure 13,  $t_{i-1}$ ,  $t_i$ , and  $t_{i+1}$  are three consecutive sampled days,  $O_{i-1}$ ,  $O_i$ ,  $O_{i+1}$  are the measured nitrate-N concentrations on these days respectively, and  $M_{i-1}$ ,  $M_i$ ,  $M_{i+1}$  are the modeled nitrate-N concentrations for these days respectively and  $R_i = O_i - M_i$  is the residual concentration on day  $t_i$ .

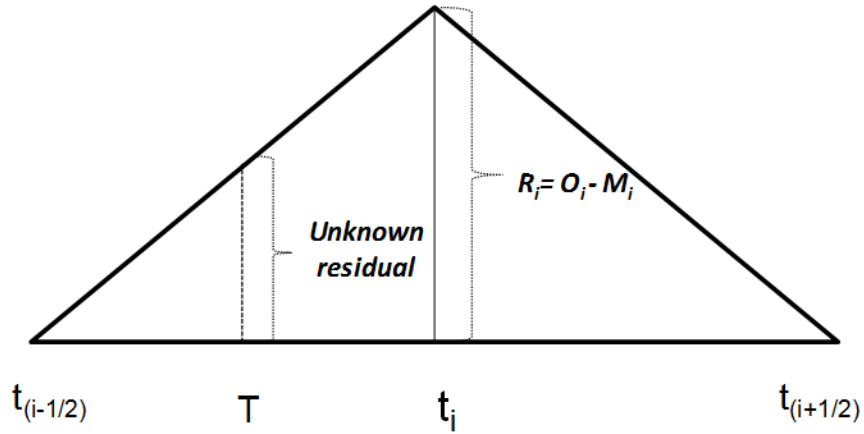


**Figure 13. Schematic of the residual distribution window.**

The residual for any time in the window,  $T$ , is calculated as follows:

- Using the Triangular distribution

Using similar triangles, Figure (14), we get the following relation:



**Figure 14. Schematic of the triangular distribution used to calculate the unknown residuals.**

$$\frac{\left[ T - t_{(i-\frac{1}{2})} \right]}{\text{Unknown residual}} = \frac{\left[ t_i - t_{(i-\frac{1}{2})} \right]}{[R_i]} \quad (12a)$$

Since,

$$R_i = O_i - M_i \quad (12b)$$

The unknown residual, for  $T = t_{i-1/2}$  to  $t_i$  can be calculated as:

$$\text{Unknown residual} = \frac{\left[ T - t_{(i-\frac{1}{2})} \right]}{\left[ t_i - t_{(i-\frac{1}{2})} \right]} [O_i - M_i] \quad (12c)$$

for  $T = t_i$  to  $t_{i+1/2}$

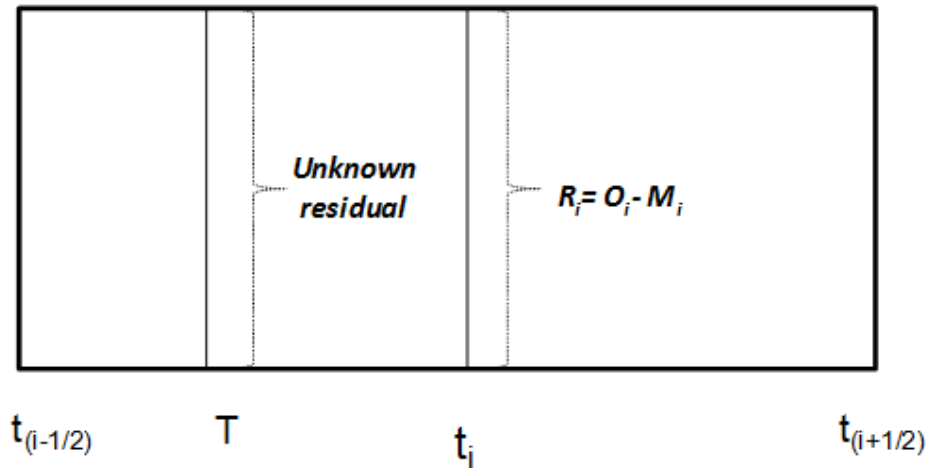
$$\text{Unknown residual} = \frac{\left[ t_{(i+\frac{1}{2})} - T \right]}{\left[ t_{(i+\frac{1}{2})} - t_i \right]} [O_i - M_i] \quad (12d)$$

- Using Rectangular distribution

For the rectangular distribution the residual remained constant throughout the window, Figure (15):

i.e. missing residuals can be calculated using equation (12b)

$$R_i = O_i - M_i$$



**Figure 15. Schematic of the rectangular distribution used to calculate the unknown residuals.**

For the proportion based error correction techniques, similar procedure was used to calculate the missing proportion values.

## 5.5 Monte Carlo Simulation

Accuracy and precision of solute load estimations are largely dependent on various parameters such as load estimation method, sampling routine, sampling frequency, storm events, and watershed size (Richards and Holloway, 1987). Uncertainties in solute load can also be attributed to human population densities in the watershed, baseflow index, and river regimes (Johnes, 2007) and can be characterized by using Monte Carlo analysis (Guo et al., 2002). In this study, sub-sampling of continuous daily concentration datasets was performed to generate sparse, replicated datasets. These sub-sampled scenarios were used to generate replicate datasets



with one-week, two-week, four-week, six-week, and eight-week sampling intervals. Monte Carlo simulation, which is essentially random number sampling, was designed based on the realization of a uniformly distributed random variate. Consequently, for any given work week, all days had an equal probability of being the sampling day. For example, if the sampling frequency was four weeks, using Monte-Carlo simulation any day within a time interval of four weeks was randomly selected as the sampling day. Also, the minimum gap between the next sampled day was set at a minimum of at least three weeks. Using this sub-sampled data, the total annual nitrate-N load was estimated by different estimators, and their precision and accuracy were evaluated against “true” or actual loads, calculated from the original datasets.

## 5.6 Evaluation Criteria

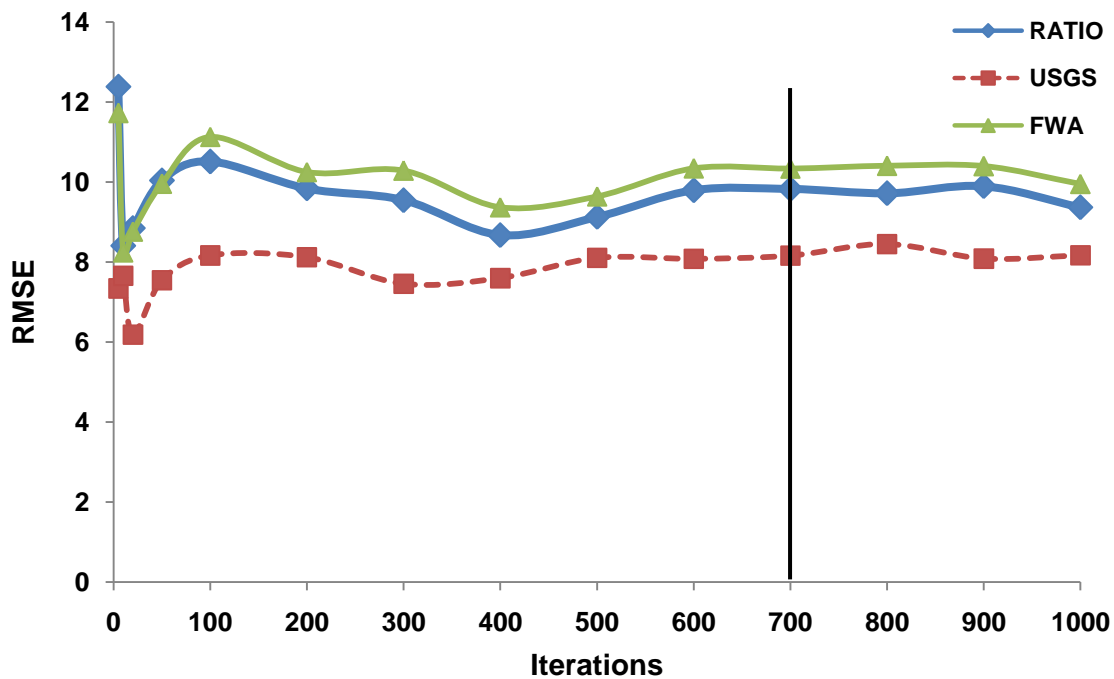
For an accurate assessment of performance of various statistical load estimation methods and error correction techniques, determination of “true” loads from contributing watersheds are of utmost importance. In this study, true loads were calculated for both watersheds using existing continuous nitrate-N and flow data from May 1, 1993 to April 30, 1999. Using daily true loads, annual loads were computed by summing the daily loads over a year. Similarly, true average yield values for both watersheds were also computed for shorter periods of records, namely, one year, two years, and three years. To evaluate the performance of load estimators and error correction techniques, bias and standard errors for load estimates were calculated using these methods. Bias, defined as the difference in estimates using an estimator and true values, was used to assess the accuracy of estimators. To gauge the precision of estimators, standard errors were computed by comparing estimates using estimators and true values. Both bias and standard errors were calculated as percentages of true loads, which permit easier comparison of estimated and true loads. As suggested by Dolan et al. (1981), Preston et al. (1989) and Guo et al. (2002), bias and standard errors were combined into an overall root mean squared estimator (RMSE) that was used to evaluate the performance of various load estimators and error correction techniques over a wide range of monitoring durations and sampling frequencies.

$$\text{RMSE} = \sqrt{B^2 + S_e^2} \quad (13)$$

To determine whether loads estimated with the different estimators converged towards the true loads on increasing or decreasing sampling frequencies, different sampling frequencies ranging from weekly to once every eight weeks were used. Monitoring durations were also varied to gauge the affect of time during which a site is monitored. Finally, RMSE values were also computed for the various error correction techniques when applied to all load estimation methods across all combinations of monitoring durations and sampling frequencies. A total of 700 iterations (Monte Carlo analysis) were performed for each simulation, and the bias and RMSE values were computed by averaging these 700 results. This was done to ensure that a wide range of observed data was used as model inputs for the three statistical load estimation methods (seven-parameter regression, ratio, and flow-weighted average). This was instrumental in better understanding of the performance of the error correction techniques, when applied to a comprehensive set of combinations of load estimator, sampling frequency, and monitoring duration.

### **5.6.1 Number of Iterations**

To eliminate any uncertainty associated with the actual date of sampling, for all monitoring scenarios 700 Monte Carlo simulations were performed and the RMSE were computed by averaging these results. This number of simulations was selected to be consistent with similar work reported by Guo et al. (2002) to facilitate a better comparison and validation of their results. Since they did not provide a rationale for using 700 simulations, the relationship between the number of simulations and the RMSE value was investigated using six years of monitoring at an eight-week sampling frequency. The result of this investigation is shown in Figure. This graph clearly indicates that after oscillating at lower values, the RMSE levels off after approximately 600 iterations, justifying the use of 700 iterations to compute the RMSE throughout the study.



**Figure 16. Number of iterations versus the RMSE for Ratio, USGS 7-parameter regression and Flow Weighted average load estimators for a scenario having six years of monitoring duration and eight week sampling frequency.**

### 5.6.2 Uncertainty Assessment

The RMSE is an appropriate indicator of the model superiority only if it can be replicated but with little variation for different realizations of the same scenario. Typically, the width of the confidence interval is proportional to the RMSE. So, when comparing two models using the RMSE it is very important to assess the changes (increase or decrease) in the width of confidence intervals or simply the absolute percentage difference in between two RMSE values. A RMSE value of 4.50% is necessarily not significantly better than a RMSE value of 5.00%, *per se*; it largely depends on the underlying marginal distributions for the RMSE values.

In the context of this study, the RMSE was used as an object function to discriminate between the performances of various models. To be certain that the RMSE statistic is sufficiently sensitive to perform this task, the variability of this statistic was determined. Two monitoring scenarios (6year-1week and 6year-8weeks) were selected and annual nitrate-N loads were calculated using the ratio estimator and the rectangular proportional error correction technique

for both the watersheds. In this study as 700 iterations were performed for each scenario to compute average bias and consequently the RMSE, 21,000 iterations were performed only for these selected scenarios to generate a larger sample size (30) of RMSE values for statistical analysis.

Table 1 summarizes the average and variance amongst these sets of RMSE for different scenarios. For these scenarios the variance was insignificant in comparison to the respective average values. Consequently, the coefficient of variance (CV) for these scenarios was also substantially low. Thus, it can be said with some certainty that the range of particular RMSE value for a 95% confidence interval is very small and thus a RMSE value of 4.00% is better than a RMSE value of 5.00%, i.e. the lower the RMSE for any load estimation method/ error correction technique/ monitoring scenario the better the technique performs. Statistical results from a few selected scenarios can be applied to all the other scenarios as all RMSE values were used to compare models whose error had the same units.

**Table 1. Uncertainty assessment of RMSE (%) values.**

<i>Vermilion River at Pontiac</i>		
	<b>6year-1week</b>	<b>6year-8weeks</b>
<b>average</b>	3.00	4.94
<b>St.Dev</b>	0.07	0.15
<b>C.V.</b>	2.33%	3.04%
<i>Upper Sangamon River at Monticello</i>		
	<b>6year-1week</b>	<b>6year-8weeks</b>
<b>average</b>	1.68	5.99
<b>St.Dev</b>	0.04	0.11
<b>C.V.</b>	2.38%	1.83%

# CHAPTER 6

## RESULTS AND DISCUSSION

### 6.1 Load Calculations

Daily flow rate and nitrate-N data from May 1, 1993 through April 30, 1999 for both watersheds were used to calculate true loads in this study. To be consistent, the analyses were limited to data from this period, as it was the only period when continuous data from both watersheds were available. The complete monitoring duration was divided into six constituent observation years (i.e., six sub-datasets of a period one year each). Using the observed datasets, annual true loads were calculated for these individual years. Total load per year for the complete monitoring duration was computed by averaging the six annual loads. Similarly, the monitoring duration was divided into three two-year periods (1993-1995, 1995-1997, 1997-1999) and two three-year periods (1993-1996, 1996-1999), and total load was calculated by averaging these sets of individual loads. Further, average yields of nitrate-N from the watersheds for the six-year monitoring durations were computed to be 26.4 kilograms per acre per year (kg/acre/year) and 23.4 kg/acre/year for Upper Sangamon at Monticello and Vermilion River at Pontiac, respectively.

### 6.2 Performance of Load Estimation Methods

RMSE values were computed for all the different load estimation methods for all possible combinations of sampling frequencies and monitoring durations. Table 2 summarizes the results obtained for Upper Sangamon at Monticello, with each RMSE value calculated by averaging results from 700 iterations. Similarly, results from Vermilion River at Pontiac are summarized in Table 3. In both tables, each box represents a unique combination of sampling frequency and monitoring duration. In each box, the first row (original) represents RMSE errors expressed as a percentage of true loads calculated using different load estimation methods, without the application of any error correction technique.

**Table 2. Summary of RMSEs given by different load estimation models and error correction techniques for different combinations of monitoring durations and sampling frequencies for the Sangamon River at Monticello. The highlighted value indicates the lowest RMSE for that particular monitoring scenario amongst all possible combinations of load estimation models and error correction techniques.**

		Sampling Frequencies															
		1 Week			2 Weeks			4 Weeks			6 Weeks			8 Weeks			
		USGS	RATIO	FLOW	USGS	RATIO	FLOW	USGS	RATIO	FLOW	USGS	RATIO	FLOW	USGS	RATIO	FLOW	
<b>Monitoring Duration</b>	<b>1 Year</b>	<b>Original</b>	4.17	3.12	3.08	7.13	4.65	4.56	27.90	7.39	7.54	NR*	8.93	9.61	NR*	NR*	NR*
		<b>Trian. Proportional</b>	3.52	3.26	3.24	6.56	4.43	<b>4.42</b>	28.96	<b>6.08</b>	6.19	NR*	<b>8.95</b>	9.06	NR*	NR*	NR*
		<b>Rect. Proportional</b>	3.75	3.37	3.29	7.01	4.62	4.69	35.58	6.68	6.80	NR*	9.76	9.83	NR*	NR*	NR*
		<b>Trian. Residual</b>	3.97	3.19	3.15	7.39	4.55	4.47	27.91	6.57	6.73	NR*	9.43	9.62	NR*	NR*	NR*
		<b>Rect. Residual</b>	3.79	3.64	3.60	7.36	5.22	5.26	30.46	7.39	7.60	NR*	10.35	10.47	NR*	NR*	NR*
		<b>Composite</b>	4.22	<b>3.05</b>	4.15	7.39	4.62	5.99	28.50	6.80	8.71	NR*	9.00	12.74	NR*	NR*	NR*
	<b>2 Years</b>	<b>Original</b>	3.48	2.61	2.67	4.97	4.11	4.16	10.33	7.03	7.04	14.82	7.51	7.63	20.55	10.38	10.51
		<b>Trian. Proportional</b>	2.59	2.37	2.40	4.15	<b>3.54</b>	3.64	10.47	5.36	5.30	14.51	<b>5.51</b>	5.60	20.22	10.80	10.81
		<b>Rect. Proportional</b>	2.80	2.44	2.47	4.91	3.60	3.63	13.36	<b>4.93</b>	4.89	18.86	6.13	6.13	23.21	10.80	<b>10.69</b>
		<b>Trian. Residual</b>	3.12	2.57	2.63	4.86	4.01	4.07	10.08	6.69	6.64	14.16	7.12	7.23	20.92	10.89	11.00
		<b>Rect. Residual</b>	2.56	2.57	2.61	4.53	3.94	4.00	11.02	6.90	6.86	15.94	7.45	7.54	23.58	12.04	12.11
		<b>Composite</b>	<b>2.29</b>	2.57	2.43	4.01	4.06	3.93	10.21	6.86	6.50	14.93	6.94	6.81	22.64	10.72	12.38
	<b>3 Years</b>	<b>Original</b>	3.31	2.35	2.27	4.35	3.08	3.02	6.86	5.69	5.83	24.12	6.53	6.41	20.13	7.33	7.34
		<b>Trian. Proportional</b>	2.41	2.04	1.96	3.41	2.66	2.58	6.39	4.33	4.45	25.10	5.68	5.65	20.66	7.14	7.13
		<b>Rect. Proportional</b>	2.71	1.89	<b>1.84</b>	4.11	2.05	<b>1.99</b>	7.73	3.50	<b>3.57</b>	34.54	<b>5.34</b>	5.37	24.55	6.03	<b>5.92</b>
		<b>Trian. Residual</b>	3.07	2.29	2.20	4.19	3.02	2.95	6.61	5.45	5.58	24.93	6.49	6.35	20.27	8.02	8.03
		<b>Rect. Residual</b>	2.64	2.17	2.11	3.96	3.00	2.88	6.61	5.67	5.79	26.64	6.63	6.58	21.72	8.74	8.73
		<b>Composite</b>	2.19	2.22	2.04	3.23	2.98	2.85	5.75	5.61	5.61	25.66	6.35	7.09	19.50	7.90	9.05
	<b>6 Years</b>	<b>Original</b>	2.58	1.96 <sup>#</sup>	1.98	4.72	3.71	3.77	5.14	4.67	4.70	11.07	5.00	5.15	7.90	7.22	6.92
		<b>Trian. Proportional</b>	1.64	1.79	1.75	3.68	3.11	3.16	<b>4.18</b>	5.27	5.25	10.39	5.98	6.06	7.64	8.73	8.58
<b>Rect. Proportional</b>		2.40	<b>1.69</b>	1.72	3.78	<b>2.01</b>	2.02	5.67	4.44	4.37	16.27	<b>3.42</b>	3.51	<b>4.76</b>	5.97	5.94	
<b>Trian. Residual</b>		2.29	1.93	1.96	4.40	3.72	3.78	4.67	4.71	4.73	10.63	5.43	5.56	7.37	7.00	6.62	
<b>Rect. Residual</b>		1.95	1.82	1.79	3.81	2.97	3.04	4.50	5.31	5.31	10.67	6.16	6.26	6.83	8.89	8.62	
<b>Composite</b>		1.75	1.85	1.73	3.70	3.52	3.43	4.11	5.02	5.21	9.66	5.47	6.40	7.96	7.88	8.79	

\*NR: No Result

Table 3. Summary of RMSEs given by different load estimation models and error correction techniques for different combinations of monitoring durations and sampling frequencies for the Vermilion River at Pontiac. The highlighted value indicates the lowest RMSE for that particular monitoring scenario amongst all possible combinations of load estimation models and error correction techniques.

		Sampling Frequencies															
		1 Week			2 Weeks			4 Weeks			6 Weeks			8 Weeks			
		USGS	RATIO	FLOW	USGS	RATIO	FLOW	USGS	RATIO	FLOW	USGS	RATIO	FLOW	USGS	RATIO	FLOW	
Monitoring Duration	1 Year	Original	6.61	5.96	5.79	9.29	8.45	8.47	28.81	10.82	10.89	NR*	13.77	13.70	NR*	NR*	NR*
		Trian. Proportional	5.92	5.72	5.55	8.87	7.53	7.52	29.19	10.46	10.50	NR*	13.39	<b>13.28</b>	NR*	NR*	NR*
		Rect. Proportional	5.00	4.97	<b>4.87</b>	8.95	6.99	<b>6.94</b>	31.29	<b>10.22</b>	10.23	NR*	13.43	13.34	NR*	NR*	NR*
		Trian. Residual	6.26	6.12	5.99	9.35	8.38	8.38	29.62	11.42	11.47	NR*	14.56	14.42	NR*	NR*	NR*
		Rect. Residual	5.82	6.37	6.20	9.93	8.96	8.95	29.70	11.84	11.79	NR*	15.22	15.01	NR*	NR*	NR*
		Composite	7.20	8.11	7.92	11.49	11.16	11.06	18.21	18.12	17.90	NR*	21.82	21.91	NR*	NR*	NR*
	2 Years	Original	5.56	4.32	4.13	6.35	6.11	6.07	12.41	8.98	8.93	30.13	11.09	11.45	14.06	12.93	13.47
		Trian. Proportional	4.51	4.01	3.95	5.46	5.28	5.26	12.33	7.84	7.73	28.69	7.75	8.11	13.46	11.71	12.18
		Rect. Proportional	3.23	3.40	<b>3.23</b>	5.06	<b>4.35</b>	4.47	14.04	5.89	<b>5.73</b>	30.29	<b>6.21</b>	6.53	15.78	<b>11.30</b>	11.59
		Trian. Residual	4.66	4.21	4.04	5.68	5.80	5.84	11.59	9.26	9.19	29.46	11.25	11.65	14.50	13.82	14.38
		Rect. Residual	3.84	4.42	4.28	5.75	6.47	6.40	11.81	9.39	9.29	29.97	10.45	10.65	17.24	14.12	14.63
		Composite	3.84	4.25	4.16	5.75	6.27	6.14	11.52	9.75	9.64	28.52	9.86	10.12	15.83	14.04	14.57
	3 Years	Original	6.09	4.30	4.31	6.82	5.76	5.55	8.19	8.12	8.31	10.92	11.91	12.00	33.38	13.70	14.29
		Trian. Proportional	5.45	4.45	4.43	6.14	5.26	5.04	7.95	7.34	7.52	11.01	9.96	10.10	35.88	12.32	12.80
		Rect. Proportional	3.22	3.16	<b>3.14</b>	4.91	<b>3.86</b>	3.80	7.62	<b>5.22</b>	5.24	10.29	<b>6.19</b>	6.31	31.46	<b>9.22</b>	9.36
		Trian. Residual	4.85	4.33	4.34	5.93	5.85	5.65	7.83	8.32	8.61	11.43	13.22	13.36	33.50	14.69	15.26
		Rect. Residual	4.17	4.58	4.57	5.68	6.23	6.04	7.08	8.57	8.86	10.77	13.56	13.59	32.98	14.69	15.30
		Composite	4.07	6.48	4.49	5.93	6.05	5.86	6.97	8.72	8.98	12.54	13.08	13.21	32.75	15.01	15.58
	6 Years	Original	4.63	3.20	3.19	5.04	4.99	4.90	6.20	6.37	6.08	8.01	9.01	9.25	8.00	9.75	9.70
		Trian. Proportional	4.00	2.72	2.65	4.24	3.90	3.82	5.64	5.79	5.47	7.60	6.98	6.93	7.91	7.96	7.81
Rect. Proportional		<b>2.36</b>	3.05	3.02	3.17	2.84	<b>2.67</b>	4.47	3.26	<b>3.04</b>	6.62	<b>3.47</b>	3.53	9.12	5.14	<b>4.96</b>	
Trian. Residual		3.51	3.21	3.19	3.92	4.99	4.90	5.48	6.49	6.17	8.32	9.65	9.91	7.51	9.82	9.92	
Rect. Residual		3.00	3.09	3.06	3.88	4.84	4.68	5.41	7.70	7.34	8.56	11.23	11.27	7.73	12.06	12.16	
Composite		2.75	3.11	3.05	4.05	4.85	4.70	5.41	7.08	6.78	8.69	10.11	10.22	7.87	11.04	11.04	

\*NR: No Result

The subsequent rows indicate RMSE errors obtained by applying various error correction techniques to the load estimation methods. In addition, the background color of each cell is correlated to the RMSE value in that particular cell. White background indicates the least RMSE values on the whole table, and an error greater than 15% is represented by a black background. All RMSE values between the least value and 15% have a scaled background color, which becomes darker with increasing magnitude and vice-versa. Any RMSE value of greater than 40% was not considered and is denoted by No Result (NR).

For both watersheds, the seven-parameter USGS method yielded appreciably high RMSE (>40%, i.e., NR) for two scenarios, namely one year of monitoring with six and eight weeks sampling frequencies. Similarly, the ratio and flow-weighted average estimator also yielded RMSE values greater than 40% for one year of monitoring and an eight weeks sampling frequency. The results indicated that accuracies of load estimation using all the three load estimators can be improved by either increasing the monitoring duration or by increasing the sampling frequency. For example, the Vermilion River at Pontiac RMSE for estimations from the seven-parameter USGS method for a sampling frequency of one week decreased from 6.61% to 4.63% on increasing the monitoring durations from one year to six years. Also, for a fixed monitoring duration of six years, increasing the sampling frequency from eight weeks to one week, the RMSE value decreased from 8.00% to 4.63%. The most accurate estimations for both watersheds were obtained with the longest monitoring duration (six years) and the highest sampling frequency (one week).

Table 4 summarizes the performance of the three load estimation methods for both watersheds used in this study. For the Vermilion River at Pontiac, the mean RMSE values for all scenarios excluding the NR scenarios were the least for the ratio estimator (8.40%). The mean RMSE for the flow-weighted average estimator was 8.45%, which is extremely close to the value for the ratio estimator. The seven-parameter USGS method yielded a mean RMSE of 11.69%. Similarly, for the Sangamon River at Monticello, the ratio estimator yielded the lowest mean RMSE for all scenarios excluding the NR scenarios at 5.44%. The flow-weighted average estimator yielded a RMSE of 5.48%, which is again very close to the mean RMSE yielded by the ratio estimator. For this watershed the seven-parameter USGS method yielded a mean RMSE of 10.20%. The minimum RMSE for both watersheds was given by the ratio estimator at 3.20%



and 1.96% for Vermilion River at Pontiac and Upper Sangamon watershed at Monticello, respectively. Both these results were obtained for the same scenario (a six-year monitoring duration and a one-week sampling frequency). The maximum RMSE for the Vermilion River at Pontiac was given by the seven-parameter USGS methods at 33.38% for a monitoring duration of three years and a sampling frequency of eight weeks. Similarly for Upper Sangamon watershed at Monticello, the maximum RMSE was 27.90% given by the seven-parameter USGS method with a monitoring duration of one year and a sampling frequency of four weeks. It should be noted that these maximum RMSE values are only for scenarios in which the RMSE was less than 40%. Standard deviations of RMSE for all possible scenarios were least for the ratio estimator at 3.75% and 2.62% for Vermilion River at Pontiac and Upper Sangamon watershed at Monticello, respectively. Similarly, standard deviations of RMSE for the flow-weighted average estimator were 3.89% and 2.69% for Vermilion River at Pontiac and Upper Sangamon watershed at Monticello, respectively. The seven-parameter USGS method had relatively high-standard deviation values, 9.30% and 8.04% for Vermilion River at Pontiac and Upper Sangamon watershed at Monticello, respectively.

**Table 4. Summary of RMSE (%) values for different load estimation models.**

<i>Vermilion River at Pontiac</i>			
	<b>7-parameter USGS</b>	<b>Ratio</b>	<b>Flow-weighted average</b>
<b>min</b>	4.63	3.20	3.19
<b>max</b>	33.38	13.77	14.29
<b>average</b>	11.69	8.40	8.45
<b>St.Dev</b>	9.30	3.75	3.89
<i>Upper Sangamon River at Monticello</i>			
	<b>7-parameter USGS</b>	<b>Ratio</b>	<b>Flow-weighted average</b>
<b>min</b>	2.58	1.96	1.98
<b>max</b>	27.90	10.38	10.51
<b>average</b>	10.20	5.44	5.48
<b>St.Dev</b>	8.04	2.62	2.69

### 6.3 Performance of Error Correction Techniques

Two sets each of two new error correction techniques, namely rectangular and triangular residual and rectangular and triangular proportional, were applied to each of the three load estimation methods for all possible combinations of sampling frequency and monitoring durations. In addition, one already existing error correction method (composite method) was also applied to the load estimation methods.

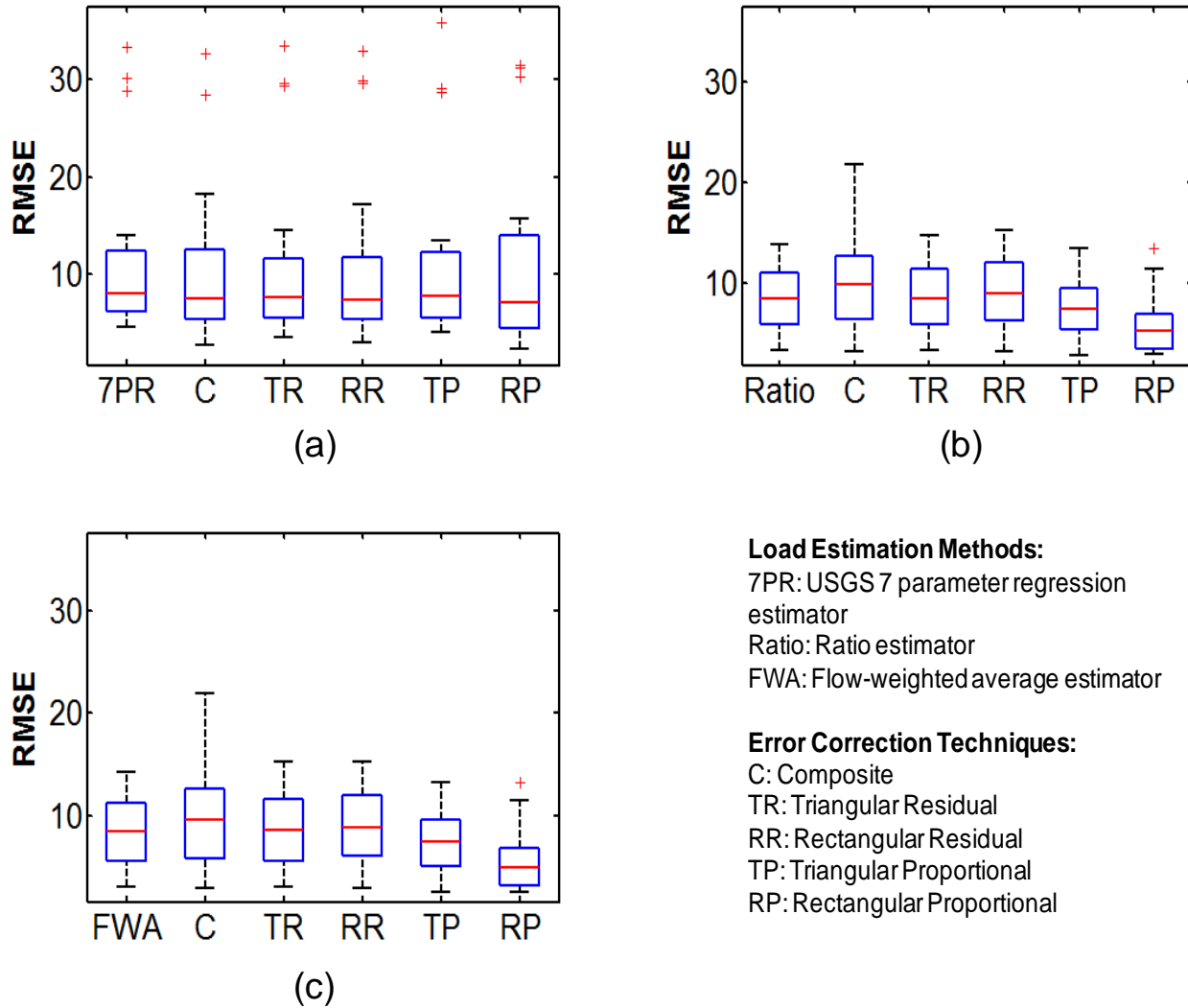
The performances of all error correction techniques applied to load estimation models are comprehensively presented in Tables 2 and 3. In general, at least one of the error correction techniques improved load estimation accuracies, except for scenarios where the seven-parameter USGS load estimation model yielded high RMSE values. For the Vermilion River at Pontiac, the least RMSE error was reduced from 3.20% to 2.36%, obtained by applying the rectangular proportional error correction technique to the seven-parameter USGS model for six-year monitoring duration with one week sampling frequency. Similarly, the lowest RMSE for the Upper Sangamon at Monticello was reduced from 1.96% to 1.64%, obtained by applying the rectangular proportional error correction technique to the seven-parameter USGS model for the same scenario. Amongst all the error correction techniques, the rectangular proportional technique, when applied to the ratio estimator, gave the least mean RMSE across all scenarios for both watersheds. The mean RMSE values for this combination were 5.90% and 4.66% for Vermilion River at Pontiac and Upper Sangamon at Monticello, respectively. These corrected results were 29.68% and 14.14% better, respectively, than the corresponding RMSE given by the uncorrected ratio estimator.

None of the new error correction techniques improve the RMSE accuracy when applied to the seven-parameter USGS method for sampling frequencies greater than two weeks in combination with monitoring durations less than six years; indeed, in some of these instances the error correction techniques increased the RMSE values. Out of 20 different possible scenarios of different monitoring durations (4) and sampling frequencies (5), only one scenario yielded RMSE values above 40% (NR) for all combinations of load estimation methods and error correction techniques applied to both watersheds. For the Vermilion River at Pontiac, the rectangular proportional error correction technique gave the most accurate estimations of annual

nitrate-N loads for 18 of the remaining 19 combinations, 9 times, 8 times and once when applied to the flow-weighted average estimator, the ratio estimator, to the seven-parameter USGS estimator, respectively. The triangular proportional error correction technique gave the least RMSE for the remaining combination. For the Upper Sangamon River at Monticello, the rectangular proportional error correction technique gave the least RMSE for nine combinations, the triangular proportional technique for five combinations, original load estimation methods without correction techniques for two combinations, and lastly, the composite method for three combinations. In general, proportional error correction techniques performed better than residual error correction techniques for both watersheds.

Box and whisker plots were created for all results from simulations for both watersheds used in this study, to visually examine the performance of various load estimation methods and error correction techniques (Figures 17 and 18). These box and whisker plots are non-parametric, with no assumptions being made about the underlying statistical distributions. Using five traditional measures; minimum, lower and upper quartile, median and maximum values, these plots provide a convenient, visual representation of the results. These box-plots (Figures 17 and 18) are indicative that the least RMSEs were obtained for the rectangular proportional error correction technique when applied to the ratio and flow-weighted average estimators. The median for this combination was the lowest for both watersheds. The box-plots also clearly indicate for the ratio and flow-weighted average load estimation techniques, the proportion based (TP, RP) error correction techniques perform better than the residual based (TR, RR, composite) error correction techniques. Interestingly, the ratio and the flow-weighted average load estimation methods gave results which were very close in magnitude with each other. This fact showed that for these watersheds, there was no appreciable improvement in load estimations by stratifying the observed data based on the flow rates as the flow-weighted average load estimator here is essentially a stratified version of the ratio load estimator.

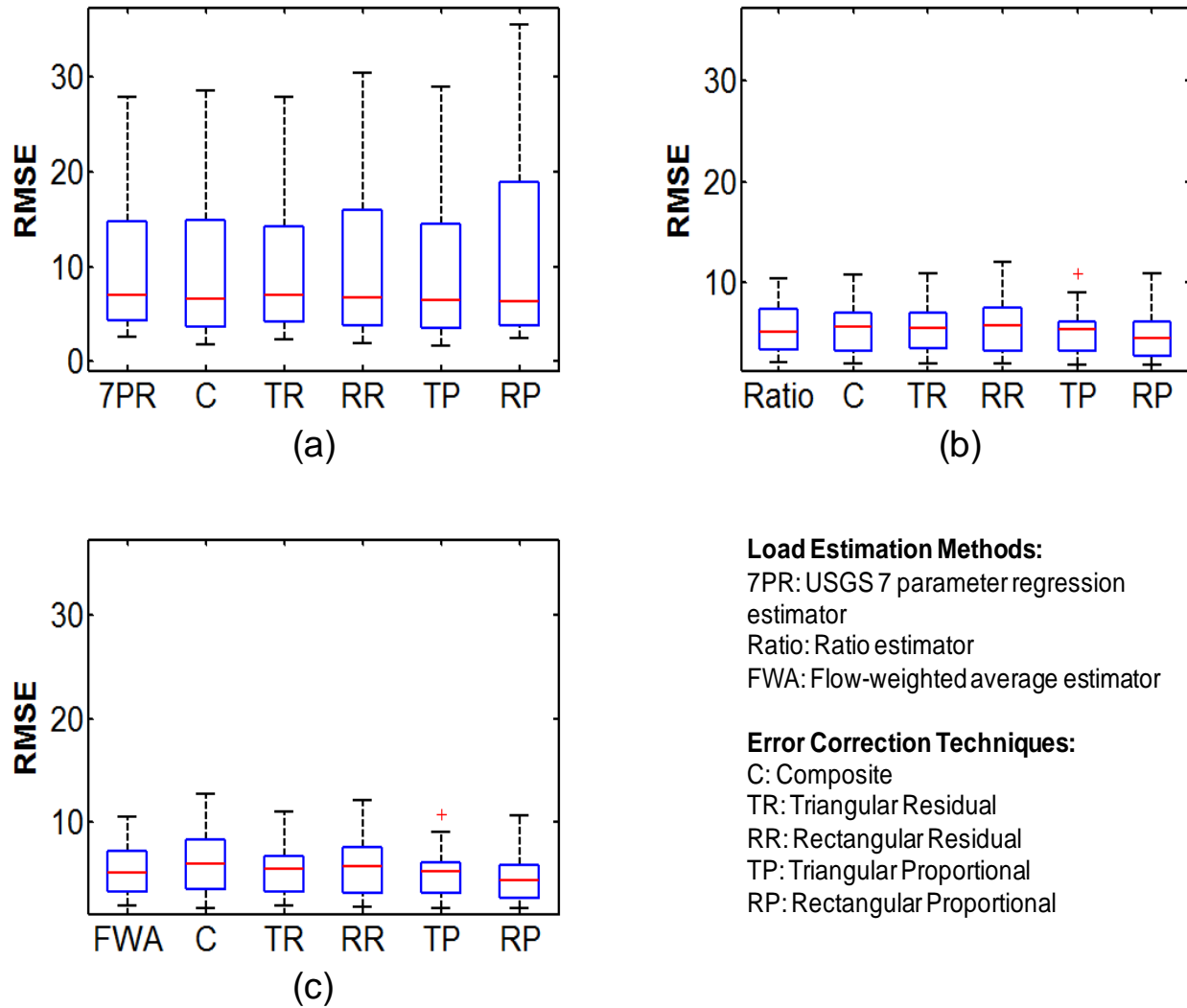
The box-plots do not give an accurate representation for the seven-parameter regression estimator. The box-plots for the seven-parameter regression estimator had a few outliers for all different conditions. For high sampling frequencies and longer monitoring durations, the seven-parameter regression estimator gives low RMSEs, and consequently the error correction techniques further improve its estimations (Tables 2 and 3).



**Figure 17. Box-plots of RMSE for different combinational scenarios of sampling frequencies (one week, two weeks, four weeks, six weeks, and eight weeks) and monitoring durations (one year, two years, three years, and six years) for different load estimation models and error correction techniques applied to them for Vermilion River at Pontiac, (a) seven-parameter regression estimator; (b) ratio estimator; (c) flow-weighted average estimator.**

For scenarios with low sampling frequencies and shorter monitoring durations, the seven-parameter regression method gives higher RMSE values; therefore the error correction techniques are not successful in improving the estimations (Tables 2 and 3). In fact, there are a few scenarios in which the error correction techniques actually increase the RMSE values in such cases. For example, for the Vermilion River at Pontiac for a monitoring duration of two years and a sampling frequency of eight weeks the seven-parameter regression method gives an RMSE

value of 14.06%, which is increased to 15.78%, on applying the rectangular proportional error correction technique (Table 3).



**Figure 18. Box-plots of RMSE for different combinational scenarios of sampling frequencies (one week, two weeks, four weeks, six weeks, and eight weeks) and monitoring durations (one year, two years, three years, and six years) for different load estimation models and error correction techniques applied to them for Sangamon River at Monticello, (a) seven-parameter regression estimator; (b) ratio estimator; (c) flow-weighted average estimator.**

## 6.4 Threshold Analysis

Evaluation of the various load estimation methods and error correction techniques in this study were conducted for four fixed sampling frequencies that did not change with flow rate, storms, or season of the year. In practice, however, sampling frequency may not be constant throughout the year, with more frequent samples taken during periods of high flow, for example. Thus, a threshold-based analysis was undertaken for both watersheds, to determine combinations of fixed and flow-threshold-based sampling that minimizes the number of samples required to achieve specified accuracy levels. Generally, sampling programs around the world are designed so that extra samples are collected during major storm events. Therefore, in such programs the total number of samples collected is a sum of samples from fixed sampling and samples from more frequent sampling when flow exceeds specified thresholds. In this study a wide range of fixed sampling frequencies were combined with flow-threshold indices ranging from 0 to 1. An index value of 0 indicates daily samples should be collected throughout the duration of a major storm, while 1 indicates that no flow-threshold sampling is done. A flow-threshold of 0.5 indicates that daily samples were also collected whenever the flow-rate was greater than or equal to 0.5 times the maximum average flow-rate for that site. Using various combinations of fixed sampling and flow-threshold sampling, sparse replicate datasets of nitrate-N concentrations were generated from continuous nitrate-N datasets by Monte Carlo sub-sampling for both watersheds. These sub-sampled datasets were generated similarly to those generated for the evaluation of load estimation methods and error correction techniques earlier, albeit in this case the sub-sampled datasets also had samples from flow-threshold sampling.

For the threshold analysis, the best performing load estimation method and error correction technique from earlier results were used to estimate nitrate-N loads for a period of six years from May 1, 1993 through April 30, 1999. The ratio estimator, along with the rectangular proportional error correction technique, was therefore used to estimate loads after the incorporation of flow-threshold sampling. As described earlier, these estimated loads were compared to the true loads and RMSE values were computed for different combinations of fixed and flow-threshold sampling. The total number of samples was also determined for each sampling combination. Figure 19(a) shows a contour plot of RMSEs for different combinations of fixed and flow-threshold sampling for Upper Sangamon River at Monticello. On the same

axes a similar contour plot was made for the number of samples, shown in Figure 19(b). Overlaying these contour plots was possible as they were plotted on the same axes. Figure 19(c) shows overlaid contour plots for Upper Sangamon River at Monticello. A similar set of contour plots for the Vermilion River at Pontiac are shown in Figure 20.

For a fixed monitoring duration, a desired accuracy could be achieved by either decreasing the sampling threshold (i.e., to capture more of the small storms), or by reducing the sampling interval in the fixed interval sampling method. For very large sampling intervals, the flow-threshold ought to be very low to achieve the desired accuracy, and similarly, for frequent sampling the flow-threshold can be relatively high. However, both of these extremes require a larger number of observations, as compared to the less extreme combinations of flow-threshold magnitude and sampling frequency.

It is fairly evident from the figures that the RMSE varies along lines for a fixed number of observations. For example, for 150 observations at the Upper Sangamon River at Monticello, the RMSE values varied from 2.25% to 5.50%. The shapes of the lines are indicative that, for a fixed number of observations, there is an "optimum" combination of threshold magnitude and sampling frequency that produces the highest accuracy for that number of observations. For example, for the Upper Sangamon River at Monticello with the overlaid contour plots, Figure 19(c) indicates that, for 150 observations, the least RMSE value of 2.50% was produced by the combination of a fixed sampling frequency of approximately 16 days and a flow-threshold of approximately 0.85. This minima, therefore, was the optimum combination of 150 observations, that produced the most accurate average nitrate-N loads over the six-year period. Such optimum combinations of sampling routines were estimated for a wide range of fixed number of total samples collected for both watersheds. These optimum combinations were indicated by the intersection of contours representing the number of samples with the contours representing the least RMSEs. These optimum points indicated the presence of a definite trend that was similar for both watersheds, as shown in Figures 19(c) and 20(c).

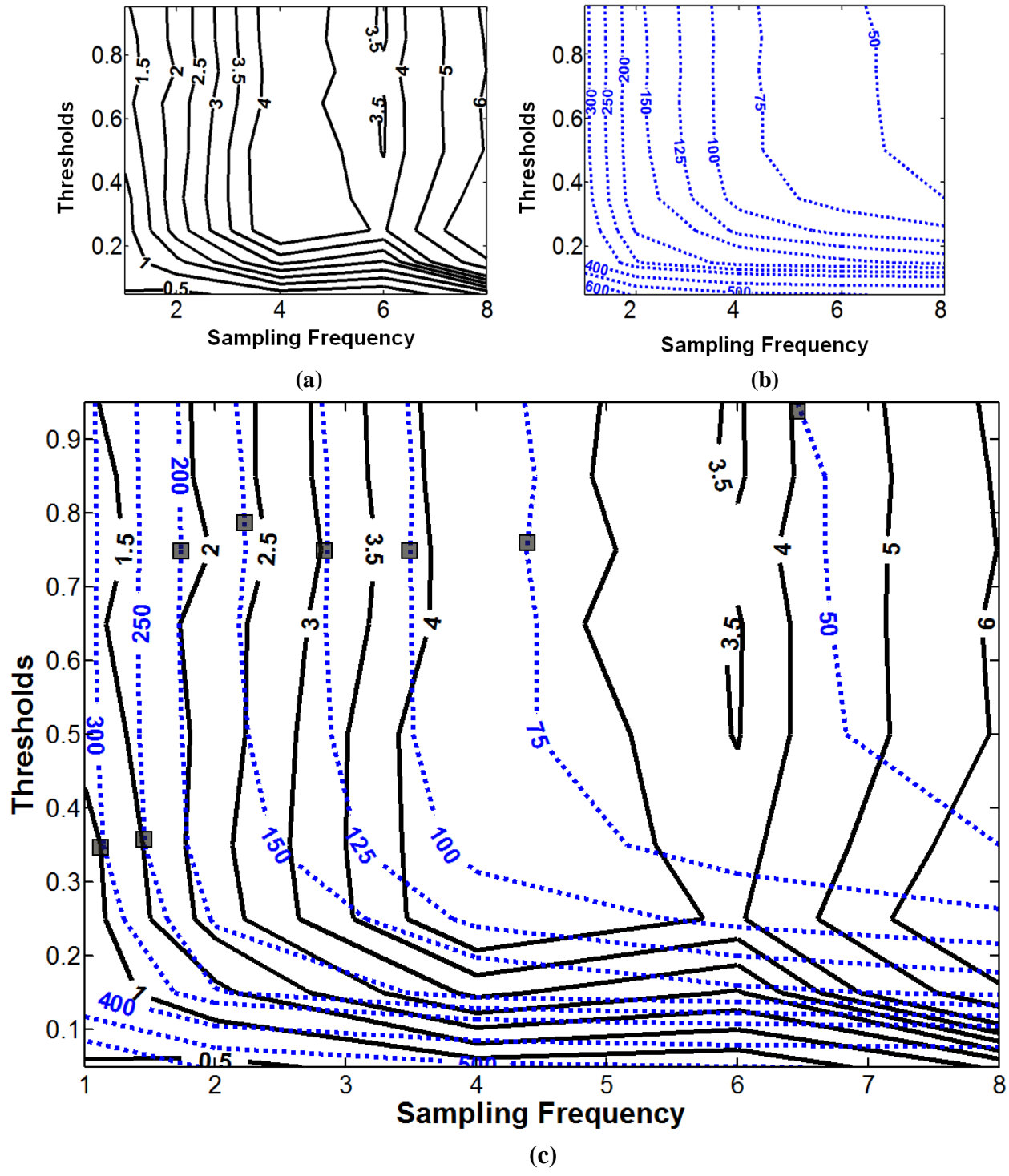


Figure 19. Contour plots for thresholds (expressed as a fraction of maximum flow over a monitoring duration from 1993-1999) and sampling frequency (weeks) of, (a) RMSE (using ratio estimator with rectangular proportional error correction technique); (b) total number of samples collected (fixed and threshold sampling); (c) super imposed contour plots of RMSE and total number of samples collected for Upper Sangamon River at Monticello.



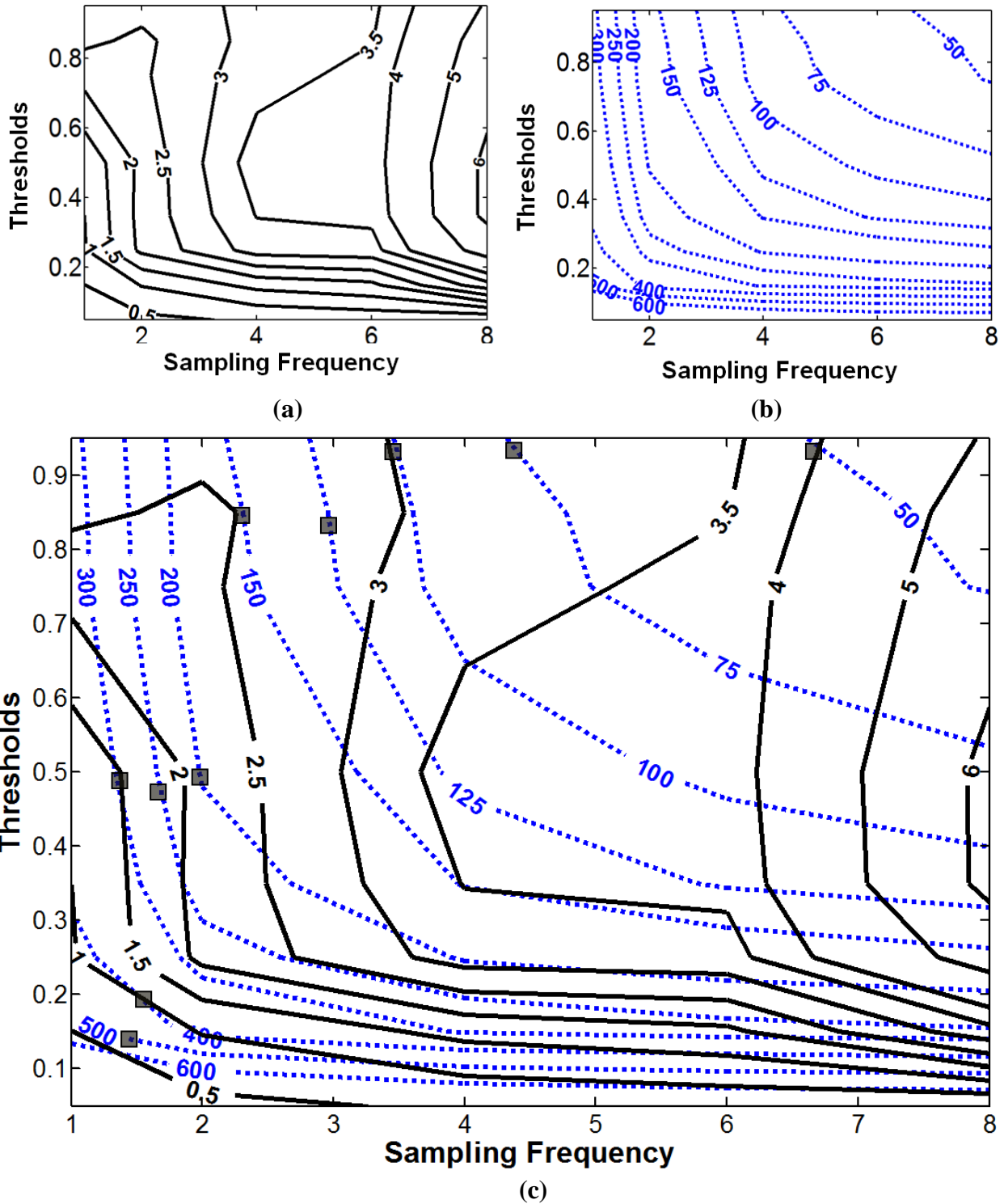


Figure 20. Contour plots for thresholds (expressed as a fraction of maximum flow over a monitoring duration from 1993-1999) and sampling frequency (weeks) of, (a) RMSE (using ratio estimator with rectangular proportional error correction technique); (b) total number of samples collected (fixed and threshold sampling); (c) super imposed contour plots of RMSE and total number of samples collected for Vermilion River at Pontiac.

# CHAPTER 7

## SUMMARY

A comprehensive analysis of the performance of three load estimation methods and five error correction techniques was conducted for two watersheds in Illinois. Sub-sampling was done from continuously recorded nitrate-N concentrations and estimated loads using various load estimation techniques, and error correction techniques were compared to true loads for different scenarios of sampling frequencies and monitoring durations.

Amongst the three load estimators, the ratio and the flow-weighted average estimators consistently gave the least RMSE values across all combinations of sampling frequencies and monitoring durations. The seven-parameter regression load estimation method gave low RMSEs for longer monitoring durations and higher sampling frequencies, but did not perform as well for scenarios with shorter monitoring durations and lower sampling frequencies. This load estimation method comprises seven different parameters, which are estimated using the sub-sampled data and then further used to estimate loads for the remaining data. In scenarios with few sub-sampled observations, it becomes very difficult to accurately calculate the seven parameters, and consequently, load estimations tend to be largely inaccurate. For example, for a two-year monitoring duration and a sampling frequency of eight weeks, only 12 sub-sampled observations were used to estimate the seven parameters. On the other hand, the ratio and flow-weighted average load estimators are not constrained by the number of parameters; they simply calculate a ratio of sub-sampled loads and use it to determine loads for the complete monitoring duration (i.e., they have a single parameter).

The application of the new error-correction techniques generally improved load estimates from the original models. Autocorrelation studies indicated that concentration residuals are auto-correlated up to a period of approximately five weeks. This fact is largely unaccounted for in the traditional statistical load estimation models. The error correction techniques presented here corrected the concentration estimates given by the models by subtracting the concentration residuals or dividing the proportion residuals depending on the technique being used. Further, because the residuals were auto-correlated, they were used to estimate residuals for unsampled

days in their vicinity, using triangular and rectangular (step) distributions. The proportional-based error correction techniques were more accurate in predicting loads, as they were better able to account for abrupt and large variations in flow rates. For the residual-based error correction techniques, the performance of the composite method was in between the triangular and rectangular residual correction methods.

The new error adjustment methods produced higher load estimation accuracies than the original unadjusted methods. To maintain the load estimation accuracy of the original methods, the new error adjustment methods will require fewer observation points than the original methods. As a consequence, the improved methodology will result in significant monitoring cost reductions. For example, for one-week sampling frequency and six-year monitoring duration, the original ratio method produced an error of 1.96%<sup>#</sup> (Table 2. Upper Sangamon River at Monticello). The error adjustment model, however, produced a similar accuracy (2.01%) for a two-week sampling frequency, reducing the number of samples by 50%.

## **CHAPTER 8**

# **CONCLUSIONS**

This study confirmed the findings of the previous research that desired load estimation accuracy could be achieved by either increasing the monitoring duration, or by sampling more frequently; the ratio and flow-weighted average methods were more accurate than the rating curve method for almost all monitoring scenarios, and that bias adjustment actually increased the biases.

This study presented new triangular and rectangular error correction techniques and applied them as proportions and residuals. The tests of the new techniques demonstrated that the best performing technique, the rectangular-proportional technique, was on average 15.13% and 30.13% relatively more accurate than the best performing original unadjusted method, the ratio estimator, for the Upper Sangamon River at Monticello and Vermilion River at Pontiac, respectively. This improved methodology can be used by the water managers to design more cost-effective monitoring plans, requiring fewer observations for the same load estimation accuracy.

Various combinations of monitoring thresholds and sampling frequencies were also tested in this study. Fixed monitoring supplemented by monitoring thresholds is the most cost-effective monitoring plan for sampling frequencies of one to three weeks. However, for lower frequencies of four or more weeks, the optimum plan is achieved without thresholds.

Relatively consistent results between the two watersheds indicated that the results could be applicable to nitrate-N load estimation at similar-sized, Midwestern streams. Nonetheless, the results of the methodology presented in this study are expected to vary with constituents, watershed size, and hydro-climatic regions. Future studies should be directed towards studying a wide range of climates, regions, constituents, and watershed sizes using the framework designed in this research.

## REFERENCES

- Alexander, R. B., R. A. Smith, G. E. Schwarz, E. W. Boyer, J. V. Nolan, and J. W. Brakebill. 2008. Differences in phosphorus and nitrogen delivery to the Gulf of Mexico from the Mississippi River Basin. *Environmental Science and Technology* 42 (3): 822-830.
- Anctil, F., M. Filion, and J. Tournebize. 2009. A neural network experiment on the simulation of daily nitrate-nitrogen and suspended sediment fluxes from a small agricultural catchment. *Ecological Modeling* 220(6): 879–887.
- Arnold, J. G., R. Srinivasan, R. S. Muttiand, and J. R. Williams. 1998. Large area hydrologic modeling and assessment part I: Model development. *Journal of the American Water Resources Association* 34(1): 73–89.
- Aulenbach, B. T., H. T. Buxton, W. A. Battaglin, and R. H. Coupe. 2007. Streamflow and nutrient fluxes of the Mississippi-Atchafalaya River Basin and subbasins for the period of record through 2005. *U.S. Geological Survey Open-File Report 2007*: 1080.
- Aulenbach, B. T., and R. P. Hooper. 2006. The composite method: An improved method for stream-water solute load estimation. *Hydrological Processes* 20(14): 3029–3047.
- Baker, D. B., R. P. Richards, T. T. Loftus, and J. W. Kramer. 2004. A new flashiness index: Characteristics and applications to Midwestern rivers and streams. *Journal of the American Water Resources Association* 40(2): 503-522.
- Beale, E. M. L. 1962. Some uses of computers in operational research. *Industrielle Organisation* 31: 51–52.
- Borah, D. K., M. Bera, and S. Shaw. 2003. Water, sediment, nutrient and pesticide measurements in an agricultural watershed in Illinois during storm events. *Transactions of the ASAE* 46(3): 657-674.

- Bramstedt, M., and T. Endres. 1999. Illinois Agricultural Experiment Station Soil Report 164. *Natural Resources Conservation Service*. Washington, D.C.
- Burn, Donald H. 1990. Real-time sampling strategies for estimating nutrient loadings. *Journal of Water Resources Planning and Management* 116 (6): 727-741.
- Byard, G. J. 2009. Effects of sampling frequency on the evaluation of nitrogen transport from subsurface drainage in Illinois. MS thesis. Urbana, IL: University of Illinois, Department of Agricultural and Biological Engineering.
- Capelo, S., F. Mira, and A. M. de Bettencourt. 2007. In situ continuous monitoring of chloride, nitrate and ammonium in a temporary stream: Comparison with standard methods. *Talanta* 71(3): 1166–1171.
- Coats, R., F. Liu, and C. R. Goldman. 2002. A Monte Carlo test of load calculation methods, Lake Tahoe Basin, California-Nevada. *Journal of the American Water Resources Association* 38(3): 719–730.
- Cochran, W. G. 1977. *Sampling Techniques* 3<sup>rd</sup> edition, John Wiley and Sons, New York.
- Cohn, T. A., D. L. Caulder, E. J. Gilroy, L. D. Zynjuk, and R. M. Summers. 1992. The validity of a simple statistical model for estimating fluvial constituent loads: An empirical study involving nutrient loads entering Chesapeake Bay. *Water Resources Research* 28(9): 2353-2563.
- Cohn, T. A., L. L. DeLong, E. J. Gilroy, R. M. Hirsch, and D. K. Wells. 1989. Estimating constituent loads. *Water Resources Research* 25(5): 937–942.
- Cookson, W. R., J. S. Rowarth, and K. C. Cameron. 2000. The effect of autumn applied 15N-labelled fertilizer on nitrate leaching in a cultivated soil during winter. *Nutrient Cycling in Agroecosystems* 56 (2): 99-107.

- David M. B., S. J. Del Grosso, X. Hu, E. P. Marshall, G. F. McIsaac, W. J. Parton, C. Tonitto, and M. A. Youssef. 2008. Modeling denitrification in a tile-drained, corn and soybean agroecosystem of Illinois, USA. *Biogeochemistry* 93 (1-2): 7-30.
- David, M. B., L. E. Gentry, D. A. Kovacic, and K. M. Smith, 1997. Nitrogen balance in and export from an agricultural watershed. *Journal of Environmental Quality* 26(4): 1038-1048.
- Dolan, D. M., A. K. Yui, and R. D. Geist. 1981. Evaluation of river load estimation methods for total phosphorus. *Journal of Great Lakes Research* 7(3): 207–214.
- Duan, N. 1983. Smearing estimator: A nonparametric retransformation method. *Journal of the American Statistical Association* 78(383): 605–610.
- Ferguson, R. I. 1986. River loads underestimated by rating curves. *Water Resources Research* 22: 74–76.
- Ferguson, R. I. 1987. Accuracy and precision of methods for estimating river loads. *Earth Surface Processes and Landforms* 12(1): 95–104.
- Goolsby, D. A., W. A. Battaglin, B. T. Aulenbach, and R. P. Hooper. 2001. Nitrogen input to the Gulf of Mexico. *Journal of Environmental Quality* 30 (2): 329–336.
- Guo, Y., M. Markus, and M. Demissie. 2002. Uncertainty of nitrate-N load computations for agricultural watersheds. *Water Resources Research* 38(10): 31–312.
- Illinois Environmental Protection Agency. 2007. *Sangamon River/ Lake Decatur Watershed TMDL Report*. IEPA.
- Illinois Environmental Protection Agency. 2009. *Vermilion River Watershed (IL Basin) TMDL Report*. IEPA.
- Im, S., K. M. Brannanand, and S. Mostaghimi. 2003. Simulating hydrologic and water quality impacts in an urbanizing watershed. *Journal of the American Water Resources Association* 9(6): 1465–1479.

- Johnes, P. J. 2007. Uncertainties in annual riverine phosphorus load estimation: Impact of load estimation methodology, sampling frequency, baseflow index and catchment population density. *Journal of Hydrology* 332(1-2): 241–258.
- Kalita, P. K., R. A. C. Cooke, S. M. Anderson, M. C. Hirschi, and J. K. Mitchell. 2007. Subsurface drainage and water quality: The Illinois experience. *Transactions of the ASABE* 50 (5):1651-1656.
- Keefer, L., and M. Demissie. 2000. Watershed Monitoring for the Lake Decatur Watershed, 1998-1999, *Illinois State Water Survey Contract Report 2000-06*, Illinois State Water Survey, Champaign, IL.
- Koch, R. W., and G. M. Smillie. 1986. Bias in hydrologic prediction using log-transformed regression models. *Water Resources Bulletin* 22(5): 717–723.
- Mukhopadhyay, B., and E. H. Smith. 2000. Comparison of statistical methods for estimation of nutrient load to surface reservoirs for sparse data set: Application with a modified model for phosphorus availability. *Water Research* 34(12): 3258–3268.
- Preston, S. D., V. J. Bierman, and S. E. Silliman. 1989. An evaluation of methods for the estimation of tributary mass loads. *Water Resources Research* 25(6): 1379-1389.
- Quilbé, R., A. N. Rousseau, M. Duchemin,, A. Poulin, G. Gangbazo, and J. P. Villeneuve. 2006. Selecting a calculation method to estimate sediment and nutrient loads in streams: Application to the Beaurivage River (Québec, Canada). *Journal of Hydrology* 326(1-4): 295–310.
- Rabalais, N. N., R. E. Turner, D. Justić, Q. Dortch, W. J. Wiseman Jr., and B. K. Sen Gupta. 1996. Nutrient changes in the Mississippi river and system responses on the adjacent continental shelf. *Estuaries* 19 (2 B): 386–407.
- Randall, G. W., W. E. Lueschen, S. D. Evans, and J. F. Moncrief. 2002. Tillage Best Management Practices for Corn-Soybean Rotations in the Minnesota River Basin, *University of Minnesota Extension*.



- Richards, R. P., and J. Holloway. 1987. Monte Carlo studies of sampling strategies for estimating tributary loads. *Water Resources Research* 23(10): 1939–1948.
- Santhi, C., J. G. Arnold, J. R. Williams, W. A. Dugas, R. Srinivasan, and L. M. Hauck. 2001. Validation of the SWAT model on a large river basin with point and nonpoint sources. *Journal of the American Water Resources Association* 37(5): 1169–1188.
- Schilling, K. E., and C. F. Wolter. 2009. Modeling nitrate-nitrogen load reduction strategies for the Des Moines River, Iowa using SWAT. *Environmental Management* 44(4): 671–682.
- Shih, G., X. Wang, H. J. Grimshaw, and J. Vanarman. 1998. Variance of load estimates derived by piecewise linear interpolation. *Journal of Environmental Engineering* 124(11): 1114–1120.
- Short, M. B. 1999. Baseline Loadings of Nitrogen, Phosphorus, and Sediments from Illinois Watersheds. *Illinois Environmental Protection Agency, IEPA/BOW/99-020, Springfield, IL.*
- Sylvan, J. B., Q. Dortch, D. M. Nelson, A. F. M. Brown, W. Morrison, and J. W. Ammerman. 2006. Phosphorus limits phytoplankton growth on the Louisiana shelf during the period of hypoxia formation. *Environmental Science and Technology* 40 (24): 7548-755.
- Turner, R. E., and N. N. Rabalais. 1994. Coastal eutrophication near the Mississippi River delta. *Nature* 368 (6472): 619–621.
- USGS. 1996. Nutrients in the Nations's Waters: Too much of a good thing? *USGS Cir 1136.* Denver CO. U.S. Geological Survey.
- Verhoff, Frank H., S. M. Yaksich, and D. A. Melfi. 1990. River nutrient and chemical transport estimation. American Society of Civil Engineers. *Journal of the Environmental Engineering Division* 106 (3): 591-608.
- Walker, W. W. 1996. Simplified Procedures for Eutrophication Assessment and Prediction: User Manual. U. S. Army Corps of Engineers, *Waterways Experiment Station, Instruction Report W-96-2.*

- Walling, D. E., and B. W. Webb. 1981. The reliability of suspended sediment load data (River Creedy, UK). Erosion and sediment transport measurement. Proc. *Florence symposium, June 1981, International Association of Hydrological Sciences, IAHS-AISH Publication 133*): 177-194.
- Wang, P., and L. C. Linker. 2008. Improvement of regression simulation in fluvial sediment loads. *Journal of Hydraulic Engineerin*, 134(10): 1527–1531.
- Young, T. C., J. V. Depinto, and T. M. Heidtke. 1988. Factors affecting the efficiency of some estimators of fluvial total phosphorus load. *Water Resources Research* 24 (9): 1535-1540.
- Yu, C., W. J. Northcott, and G. F. McIsaac. 2004. Development of AN artificial neural network for hydrologic and water quality modeling of agricultural watersheds. *Transactions of the ASAE* 47(1): 285–290.

# APPENDIX A.

## A.1. Matlab code to calculate bias and RMSE for various load estimation models and error correction techniques using 700 iterations.

### Part 1

```
% Before running the program YOU MUST load the files, for example:
% load N9394.txt
% load Q9394.txt
% part1(N9394,Q9394)
function
[m,Triangular_Ratio,Triangular_Residual,Square_Ratio,Square_Residual,Composite,USGS_Regression,Ratio,Flow
_Weighted]=Vmc2w(n9999,q9999)
% Monte Carlo simulation to determine the accuracy (bias) and precision
% (standard deviation) of various load estimation methods with different
% sampling frequencies.
true9394 = 43.792; true9495 =20.022; true9596 = 21.225; true9395 = 31.907;
true9697 = 38.362; true9798 =24.469; true9899 = 37.138; true9597 = 29.793;
true9399 = 30.835; true9799 = 30.804; true9396 = 28.346; true9699 = 33.323;
if inputname(1)=='n9699'
    true = true9699
end
if inputname(1)=='n9396'
    true = true9396
end
if inputname(1)=='n9395'
    true = true9395
end
if inputname(1)=='n9597'
    true = true9597
end
if inputname(1)=='n9799'
    true = true9799
end
if inputname(1)=='n9394'
    true = true9394
end
if inputname(1)=='n9495'
    true = true9495
end
if inputname(1)=='n9596'
    true = true9596
end
if inputname(1)=='n9697'
    true = true9697
end
if inputname(1)=='n9798'
    true = true9798
```

```

end
if inputname(1)=='n9899'
    true = true9899
end
if inputname(1)=='n9399'
    true = true9399
end
for n = 1:700
    n

[TRAT_1(n),SRAT_1(n),TRES_1(n),SRES_1(n),comp_1(n),regr_1(n),ratio_1(n),flow_w_1(n)]=part2(n9999,q9999);
end
% meanD = mean(Demissie_1);
meanTrat = mean(TRAT_1);
meanSrat = mean(SRAT_1);
meanTres = mean(TRES_1);
meanSres = mean(SRES_1);

meanP = mean(comp_1);
meanC = mean(regr_1);

meanR = mean(ratio_1);
meanF = mean(flow_w_1);
% stdD = std(Demissie_1);
% biasD = meanD - true;
% rmseD = sqrt(biasD^2 + stdD^2);
stdTrat = std(TRAT_1);
biasTrat = meanTrat - true;
rmseTrat = sqrt(biasTrat^2 + stdTrat^2);

stdSrat = std(SRAT_1);
biasSrat = meanSrat - true;
rmseSrat = sqrt(biasSrat^2 + stdSrat^2);

stdTres = std(TRES_1);
biasTres = meanTres - true;
rmseTres = sqrt(biasTres^2 + stdTres^2);

stdSres = std(SRES_1);
biasSres = meanSres - true;
rmseSres = sqrt(biasSres^2 + stdSres^2);

stdP = std(comp_1);
biasP = meanP - true;
rmseP = sqrt(biasP^2 + stdP^2);

stdC = std(regr_1);
biasC = meanC - true;
rmseC = sqrt(biasC^2 + stdC^2);

stdR=std(ratio_1);
biasR = meanR - true;
rmseR = sqrt(biasR^2 + stdR^2);

```

```

stdF = std(flow_w_1);
biasF = meanF - true;
rmseF = sqrt(biasF^2 + stdF^2);
% Demissie_STD_bias_rmse =[stdD,biasD,rmseD]*100/true
Triangular_Ratio = [stdTrat,biasTrat,rmseTrat]*100/true
Square_Ratio = [stdSrat,biasSrat,rmseSrat]*100/true
Triangular_Residual = [stdTres,biasTres,rmseTres]*100/true
Square_Residual = [stdSres,biasSres,rmseSres]*100/true

Composite = [stdP,biasP,rmseP]*100/true
USGS_Regression = [stdC,biasC,rmseC]*100/true
Ratio =[stdR,biasR,rmseR]*100/true
Flow_Weighted =[stdF,biasF,rmseF]*100/true

SD=[stdTrat,abs(biasTrat),rmseTrat; stdSrat,abs(biasSrat),rmseSrat; stdTres,abs(biasTres),rmseTres;
    stdSres,abs(biasSres),rmseSres; stdP,abs(biasP),rmseP; stdC,abs(biasC),rmseC; stdR,abs(biasR),rmseR;
    stdF,abs(biasF),rmseF]*100/true;
bar(SD,'group')
set(gca,'XTickLabel',{'Trian_Ratio','Sq_Ratio','Trian_Residual','Sq_Residual','Composite','USGS','Ratio','Flow'})
legend('St.Deviation','Bias','RMSE',2)
colormap summer;

save testdata1 SRAT_1 TRAT_1
%MM AFTER EACH RUN THE SCREEN WILL DISPLAY THE FIRST VARIABLE, IN THIS CASE
C_STD_b_r
cputime
tic
toc
clear

```

## A.2. Matlab code to compute load estimates from various load estimation models and error correction techniques.

### Part 2

```
function [TRAT_1,SRAT_1,TRES_1,SRES_1,comp_1,regr_1,ratio_1,flow_w_1] = part2(n9999,q9999)
% Monte Carlo simulation program for weekly sampling frequency. Annual
% average load is calculated using two averaging methods, one ratio method
% and two regression methods. The two regression methods are both based on
% Cohn's seven parameter regression equation with one of them corrected for
% bias based on the MVUE suggested by Cohn and the other corrected for bias
% based on the smearing estimator.
% The units used in the input files are cfs and mg/L, One cubic foot = 28.32 liters.
global ymvue ysmear
constant = 28.32*3600*24/1.e+9;
% The Watershed area is 370,560 acres. One Mg = 2205 pounds
area = 370560; Mg2Pounds = 2205;
% week is an external function.
w = week(2,n9999);
R = ceil(5*rand(318,1));
% sw_n - selected week number
% nd - number of day in the record
% nw - number of weeks in the record
% ny - number of years in the record
sw_n = 0;
sn = 1;
nd = size(w,1);
nw = w(nd,2);
ny = nd/365;
bw = ceil(2*rand(1,1));
% Randomly select a day (Monday through Friday) and store NO3 concentration
% and discharge data for that day in array s.
for wn = 1:nw
    tn = 0;
    for n = sn:nd
        if w(n,2) == wn
            tn = tn + 1;
            temp(tn,1) = w(n,1);
            temp(tn,2) = w(n,2);
            temp(tn,3) = w(n,3);
        end
    end
    if wn == bw
        for n = 1:tn
            if temp(n,3) == R(wn)
                sw_n = 1;
                s(sw_n,1) = temp(n,1);
                s(sw_n,2) = wn;
                s(sw_n,3) = n9999(temp(n,1),2);
                s(sw_n,4) = q9999(temp(n,1),2);
                s(sw_n,5) = (s(sw_n,4)*s(sw_n,3))*(constant)/(ny*area)*Mg2Pounds;
            end
        end
    end
end
```

```

        end
    end
end
g=1;
if rem(wn-bw,g) == 0 && wn > bw
for n = 1:tn
    if temp(n,3) == R(wn)
        swn = swn + 1;
        s(swn,1) = temp(n,1);
        s(swn,2) = wn;
        s(swn,3) = n9999(temp(n,1),2);
        s(swn,4) = q9999(temp(n,1),2);
        s(swn,5) = (s(swn,4)*s(swn,3))*(constant)/(ny*area)*Mg2Pounds;

        end
    end
end

sn = sn + tn;
end

% Prepare to calculate load using Cohn's regression method
q_sample_tot = 0;
l_sample_tot = 0;
for n = 1:swn
    q_sample_tot = q_sample_tot + s(n,4);
    l_sample_tot = l_sample_tot + s(n,3)*s(n,4);
    q(n,1) = log(s(n,4));
    t(n,1) = (s(n,1)+120)/365;
end
tmean = mean(t);
qmean = mean(q);
cubsumt = 0;
sqsumt = 0;
cubsumq = 0;
sqsumq = 0;
for n = 1:swn
    cubsumt = cubsumt + (t(n,1)-tmean)^3;
    sqsumt = sqsumt + (t(n,1)-tmean)^2;
    cubsumq = cubsumq + (q(n,1)-qmean)^3;
    sqsumq = sqsumq + (q(n,1)-qmean)^2;
end
tcenter = tmean + cubsumt/(2*sqsumt);
qcenter = qmean + cubsumq/(2*sqsumq);

for n = 1:swn
    x(n,1) = 1;
    x(n,2) = q(n,1) - qcenter;
    x(n,3) = (x(n,2))^2;
    x(n,4) = sin(2*pi*t(n,1));
    x(n,5) = cos(2*pi*t(n,1));
    y(n,1) = log(s(n,3));
end
[b,bint,r,rint,stats] = regress(y,x);

```

```

std_of_r=std(r);

% after obtaining coefficients of regression, determine NO3
% concentration using the obtained regression equation

for n = 1:nd
    q(n,1) = log(q9999(n,2));
    t(n,1) = (q9999(n,1)+120)/365;
end
tmean = mean(t);
qmean = mean(q);
cubsumt = 0;
sqsumt = 0;
cubsumq = 0;
sqsumq = 0;
for n = 1:nd
    cubsumt = cubsumt + (t(n,1)-tmean)^3;
    sqsumt = sqsumt + (t(n,1)-tmean)^2;
    cubsumq = cubsumq + (q(n,1)-qmean)^3;
    sqsumq = sqsumq + (q(n,1)-qmean)^2;
end
tcenter = tmean + cubsumt/(2*sqsumt);
qcenter = qmean + cubsumq/(2*sqsumq);

for n = 1:nd
    x(n,1) = 1;
    x(n,2) = q(n,1) - qcenter;
    x(n,3) = (x(n,2))^2;
    x(n,4) = sin(2*pi*t(n,1));
    x(n,5) = cos(2*pi*t(n,1));
end
y = exp(x*b);

% stats
% Determine bias correction factors of MVUE and Smearing estimator
v = x*(inv(x'*x))*(x');
s2 = (swn-1)*var(r)/(swn-7);
c = exp(r);
SM = sum(c)/swn;
m=swn-5;
for n = 1:nd
    arg = (m+1)*(1-v(n,n))*s2/(2*m);
    Gm(n) = gm(m,arg);
    ymvue(n) = y(n)*Gm(n);
    ysmear(n) = y(n)*SM;
end

% Calculate loads using regression and ratio methods.
q_tot = 0;
% trueload = 0;
regrload = 0;
mvueload = 0;
smearload = 0;
composite= 0;

```



```

residuals=0;
B=0;C=0;
B=zeros(length(n9999),2);
MB1=zeros(length(n9999),2);
MB2=zeros(length(n9999),2);
MB3=zeros(length(n9999),2);
MB4=zeros(length(n9999),2);
for n=1:nd
    B(n,1)=n;
    B(n,2)=(constant/(ny*area)*Mg2Pounds*q9999(n,2)*y(n));
    MB1(n,1)=n;
    MB1(n,2)=(constant/(ny*area)*Mg2Pounds*q9999(n,2)*y(n));
    MB2(n,1)=n;
    MB2(n,2)=(constant/(ny*area)*Mg2Pounds*q9999(n,2)*y(n));
    MB3(n,1)=n;
    MB3(n,2)=(constant/(ny*area)*Mg2Pounds*q9999(n,2)*y(n));
    MB4(n,1)=n;
    MB4(n,2)=(constant/(ny*area)*Mg2Pounds*q9999(n,2)*y(n));
end

sb=size(MB1);

sz=size(s);

R1=zeros(length(n9999),2);
R2=zeros(length(n9999),2);
R3=zeros(length(n9999),2);
R4=zeros(length(n9999),2);
for n = 1:nd
    R1(n,1) = q9999(n,1);
    R1(n,2) = (l_sample_tot/q_sample_tot)*q9999(n,2)*constant/(ny*area)*Mg2Pounds;
    R2(n,1) = q9999(n,1);
    R2(n,2) = (l_sample_tot/q_sample_tot)*q9999(n,2)*constant/(ny*area)*Mg2Pounds;
    R3(n,1) = q9999(n,1);
    R3(n,2) = (l_sample_tot/q_sample_tot)*q9999(n,2)*constant/(ny*area)*Mg2Pounds;
    R4(n,1) = q9999(n,1);
    R4(n,2) = (l_sample_tot/q_sample_tot)*q9999(n,2)*constant/(ny*area)*Mg2Pounds;
end
C=zeros(sz(1,1),2);
CC=zeros(sz(1,1),2);
for i=1:sz(1,1)
    for j=1:length(R1)
        if s(i,1)==R1(j,1)
            C(i,1)=s(i,1);
            C(i,2)=R1(j,2)/s(i,5);
            CC(i,1)=s(i,1);
            CC(i,2)=R1(j,2)-s(i,5);
        end
    end
end
end

sc=size(C);

P=[];

```

```
SRAT=[];
```

```
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%  
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%% SQUARE RATIO METHOD
```

```
t1=3*g;  
for i=2:sc(1,1)-1  
    for j=1:sb(1,1)  
        if C(i,1)==R1(j,1)  
            for k=-t1:t1  
                R1(j+k,2)=R1(j+k,2)/C(i,2);  
            end  
        end  
    end  
end  
SRAT_1=sum(R1(:,2));
```

```
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%  
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
```

```
K=[];  
TRAT=[];
```

```
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%  
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%% TRIANGULAR RATIO METHOD
```

```
K=zeros(5000,2);  
t2=3*g;  
for i=2:sc(1,1)  
    X=[ C(i,1)-t2; C(i,1); C(i,1)+t2];  
    Y=[ 1; C(i,2); 1];  
    Xi=(C(i,1)-t2):(C(i,1)+t2);  
    Xi=Xi';  
    Yi=interp1(X,Y,Xi);  
    for j=1:length(Xi)  
        K(Xi(j),1)=Xi(j);  
        K(Xi(j),2)=Yi(j);  
    end  
end
```

```
end  
  
K1=K(:,1);  
P(:,1) = K1(K1~=0);  
K2=K(:,2);  
P(:,2) = K2(K2~=0);
```

```
for i=2:length(P)-2  
    for j=1:sb(1,1)-10  
        if P(i,1)==R2(j,1)  
            R2(j,2)=R2(j,2)/P(i,2);  
        end  
    end  
end  
end
```

```
TRAT_1=sum(R2(:,2));
```

```
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
```

```
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%  
% TRIANGULAR RESIDUAL METHOD
```

```
t3=3*g;  
KT=zeros(370*6,2);  
for i=1:370*6  
    KT(i,1)=-100000;  
    KT(i,2)=-100000;  
end  
sk=length(KT);  
for i=3:sc(1,1)  
    X=[ CC(i,1)-t3; CC(i,1); CC(i,1)+t3];  
    Y=[ 0; CC(i,2); 0];  
    Xi=(CC(i,1)-t3):(CC(i,1)+t3);  
    Xi=Xi';  
    Yi=interp1(X,Y,Xi);  
    for j=1:length(Xi)  
        if KT(Xi(j),1)==-100000  
            KT(Xi(j),1)=Xi(j);  
            KT(Xi(j),2)=Yi(j);  
        else  
            KT(Xi(j),1)=Xi(j);  
            KT(Xi(j),2)=(KT(Xi(j),2)+Yi(j))/2;  
        end  
    end  
end  
% KT;  
sK=size(KT);
```

```
for i=1:sb(1,1)  
    for j=1:sK(1,1)  
        if R3(i,1)==KT(j,1)  
            R3(i,2)= R3(i,2)-KT(j,2);  
        end  
    end  
end
```

```
R3;  
TRES_1=sum(R3(:,2));
```

```
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
```

```
SRES=[];  
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%  
% SQUARE RESIDUAL METHOD
```

```
t4=3*g;  
for i=2:sc(1,1)-1
```

```

for j=1:sb(1,1)
    if CC(i,1)==R4(j,1)
        for k=-t4:t4
            R4(j+k,2)=R4(j+k,2)-CC(i,2);
        end
    end
end
end
SRES_1=sum(R4(:,2));
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

```

```

F(1,1)=B(1,1);
F(1,2)=B(1,2);
F(sc+2,1)=B(nd,1);
F(sc+2,2)=B(nd,2);
for i=1:sc
    F(i+1,1)=CC(i,1);
    F(i+1,2)=CC(i,2);
end
for d=1:nd
    L(d,1)=interp1q(F(:,1),F(:,2),d);
end

```

```
sL=size(L);
```

```

for n = 1:nd
    q_tot = q_tot + q9999(n,2);
    % trueload = trueload + Q9999(n,2)*N9999(n,2);
    % residuals(n,1) = ((constant/(ny*area)*Mg2Pounds*q9999(n,2)*y(n))-L(n,2));
    composite(n,1) = B(n,2)-L(n,1);
    regrload = regrload + q9999(n,2)*y(n);
    mvueload = mvueload + q9999(n,2)*ymvue(n);
    smearload = smearload + q9999(n,2)*ysmear(n);
end
% true_load = trueload*constant/(ny*area)*Mg2Pounds;
composite=composite(~isnan(composite));
%ncomp_1=sum(B);
comp_1=sum(composite);
regr_1 = regrload*constant/(ny*area)*Mg2Pounds;
mvue_1 = mvueload*constant/(ny*area)*Mg2Pounds;
smear_1 = smearload*constant/(ny*area)*Mg2Pounds;
ratio_1 = (l_sample_tot/q_sample_tot)*q_tot*constant/(ny*area)*Mg2Pounds;
% Calculate load using stratified flow-weighted average concentration
% method, which is the method used by Matthew Short of IEPA based on
% the program developed by US Army Corps of Engineers (Walker, 1996).
meanq = q_tot/nd;
half_meanq = 0.5*meanq;
doub_meanq = 2*meanq;
tot_l1 = 0;
tot_l2 = 0;
tot_l3 = 0;
tot_q1 = 0;
tot_q2 = 0;

```

```

tot_q3 = 0;
for n = 1:swn
    if s(n,4) <= half_meanq
        tot_l1 = tot_l1 + s(n,4)*s(n,3);
        tot_q1 = tot_q1 + s(n,4);
    end
    if s(n,4) > half_meanq & s(n,4) <= doub_meanq
        tot_l2 = tot_l2 + s(n,4)*s(n,3);
        tot_q2 = tot_q2 + s(n,4);
    end*
    if s(n,4) > doub_meanq
        tot_l3 = tot_l3 + s(n,4)*s(n,3);
        tot_q3 = tot_q3 + s(n,4);
    end
end
mean_n2 = tot_l2/tot_q2;
if tot_q1 == 0
    mean_n1 = mean_n2;
else
    mean_n1 = tot_l1/tot_q1;
end
if tot_q3 == 0
    mean_n3 = mean_n2;
else
    mean_n3 = tot_l3/tot_q3;
end
ave_l1 = 0;
ave_l2 = 0;
ave_l3 = 0;
for n = 1:nd
    if q9999(n,2) <= half_meanq
        ave_l1 = ave_l1 + mean_n1*q9999(n,2);
    end
    if q9999(n,2) > half_meanq & q9999(n,2) <= doub_meanq
        ave_l2 = ave_l2 + mean_n2*q9999(n,2);
    end
    if q9999(n,2) > doub_meanq
        ave_l3 = ave_l3 + mean_n3*q9999(n,2);
    end
end

for n = 1:nd
    F(n,1)=n;
    if q9999(n,2) <= half_meanq
        F(n,2)= mean_n1*q9999(n,2);
    end
    if q9999(n,2) > half_meanq & q9999(n,2) <= doub_meanq
        F(n,2)= mean_n2*q9999(n,2);
    end
    if q9999(n,2) > doub_meanq
        F(n,2)= mean_n3*q9999(n,2);
    end
end
end
flow2 = (sum(F(:,2)))*constant/(ny*area)*Mg2Pounds;
flow_w_1 = (ave_l1+ave_l2+ave_l3)*constant/(ny*area)*Mg2Pounds;

```

```
for i=1:length(n9999)

    actualload(i) = (n9999(i,2)*q9999(i,2)*constant*2205)/area/ny;
    ysave(i) = (y(i)*constant*2205*q9999(i,2))/area;

end
A=actualload';
R=B(:,2);

save testdata R A L composite C MB1 MB2 MB3 MB4 s B y x b CC K P
```